

La Salle University La Salle University Digital Commons

Mathematics and Computer Science Capstones

Mathematics and Computer Science, Department
of

Summer 8-31-2015

Storage and Analysis of Big Data Tools for Sessionized Data

Robert McGinley

La Salle University, mcginleyr1@lstudent.lasalle.edu

Jason Etter

La Salle University, etterj1@student.lasalle.edu

Follow this and additional works at: <http://digitalcommons.lasalle.edu/mathcompcapstones>



Part of the [Databases and Information Systems Commons](#)

Recommended Citation

McGinley, Robert and Etter, Jason, "Storage and Analysis of Big Data Tools for Sessionized Data" (2015). *Mathematics and Computer Science Capstones*. 24.

<http://digitalcommons.lasalle.edu/mathcompcapstones/24>

This Thesis is brought to you for free and open access by the Mathematics and Computer Science, Department of at La Salle University Digital Commons. It has been accepted for inclusion in Mathematics and Computer Science Capstones by an authorized administrator of La Salle University Digital Commons. For more information, please contact careyc@lasalle.edu.

INL 880 Capstone Project: Storage and Analysis of Big Data Tools for Sessionized Data

Supervised by: Professor Peggy McCoey

By

Robert McGinley & Jason Etter

Final Draft: 7/20/2015

Executive Summary

The Oracle database currently used to mine data at PEGGY is approaching end-of-life and a new infrastructure overhaul is required. It has also been identified that a critical business requirement is the need to load and store very large historical data sets. These data sets contain raw electronic consumer events and interactions from a website such as page views, clicks, downloads, return visits, length of time spent on pages, and how they got to the site / originated. This project will be focused on finding a tool to analyze and measure sessionized data, which is a unit of measurement in web analytics that captures either a user's actions within a particular time period, or the process of segmenting user activity of each user into sessions, each representing a single visit to the site. This sessionized data can be used as the input for a variety of data mining tasks such as clustering, association rule mining, sequence mining etc (Ansari. 2011) This sessionized data must be delivered in a reorganized and readable format timely enough to make informed go-to-market decisions as it relates to the current and existing industry trends. It is also pertinent to understand any development work required and the burden on the resources.

Legacy on-premise data warehouse solutions are becoming more expensive, less efficient, less dynamic, and unscalable when compared to current Cloud Infrastructure as a Service (IaaS) that offer real time, on-demand, pay-as-you-go solutions . Therefore, this study will examine the total cost of ownership (TCO) by considering, researching, and analyzing the following factors against a system wide upgrade of the current on-premise Oracle Real Application Cluster (RAC) System:

- High performance: real-time (or as close to as possible) query speed against sessionized data
- SQL compliance
- Cloud based or, at least a hybrid (read: on-premise paired with cloud)
- Security: encryption preferred
- Cost structure: cost-effective pay-as-you-go pricing model and resources required for the migration and operations.

These technologies analyzed against the current Oracle database are:

- Amazon Redshift
- Google Bigquery
- Hadoop
- Hadoop + Hive

The cost of building an on-premise data warehouse is substantial. The project will determine the performance capabilities and affordability of Amazon Redshift, when compared to other emerging highly ranked solutions, for running e-commerce standard analytics queries on terabytes of sessionized data. Rather than redesigning, upgrading, or over purchasing infrastructure at a high cost for an on-premise data warehouse, this project considers data warehousing solutions through cloud based infrastructure as a service (IaaS) solutions. The proposed objective of this project is to determine the most cost-effective high performer between Amazon Redshift, Apache Hadoop, and Google BigQuery when running e-commerce standard analytics queries on terabytes of sessionized data.

Introduction

PEGGY, an E Commerce organization that sells vintage music memorabilia, has collected thirty (30) terabytes of data that represent the recorded sessions of the user's interactions with the PEGGY online store.

Disk Storage
· 1 Bit = Binary Digit
· 8 Bits = 1 Byte
· 1000 Bytes = 1 Kilobyte
· 1000 Kilobytes = 1 Megabyte
· 1000 Megabytes = 1 Gigabyte
· 1000 Gigabytes = 1 Terabyte
· 1000 Terabytes = 1 Petabyte

source: What's A Byte?, "Megabytes, Gigabytes, Terabytes... What are they?"

By combining the massive amounts of data captured from the PEGGY user groups, along with newer more powerful analytics and regression algorithms, there is a greater chance to predict future outcomes. There is an expectation to utilize real-time insights by automating or providing a short list of actions to significantly improve business growth. (Minelli. 2012)

The current Oracle based data warehouse serves as the core system that fetches, analyzes, and readies all of this data for business reporting. It has been identified that this traditional on-premise data warehousing system, although reliable, requires significant engineering overhead to cleanse, transform and insert for later use in data aggregation and analysis. In many cases, the data reports are virtually unusable when compared to competing organizations capable of bulk loading and cleansing data automatically through cloud based solutions.

To note, updating a new system, such as a traditional data warehouse could take several years along with a significant investment in resources to be configured and completed. Additionally, managing and administering the current data warehouse requires significant time and resources. Over 70% of the technology budget spent for on-premise systems are drained before the system is fully functional. (Lohr. 2012) This includes the hardware and software to be installed, multiple components required to be optimized and customized, and updating and maintaining current technology. Further, it is estimated that the upfront costs of the database will cost approximately \$4,850,000.00. (MongoDB. 2015).

In contrast, cloud solutions like Redshift range between \$65,000.00 and \$132,000.00 per year for comparable infrastructure. It is imperative to compare the current data warehousing systems against the newly available cloud and open source solutions. In addition to upfront costs, this paper will assess the cost of training and resource allocation required for these tools to ensure the total cost of ownership (TCO) matches adequately. Lastly, it has been identified that the current querying speeds of the Oracle database compared to the potential of real-time data analytics in the cloud could provide user traffic reporting at a speed that would currently be inconceivable.

The Value of Sessionized Data

The information learned from recording and analyzing how people browse online are known as web logs which are used to reconstruct the path on any given website. A stream of these web logs, more widely known as sessions, is a stream of records regarding the user's individual clicks. These clicks are known as the clickstream, which is then sessionized resulting

in a vector. A vector is a quantifiable magnitude headed in a specific direction. In this case, the clickstream results in a vector and can be recorded, analyzed, and compared. (Liebowitz)

To serve as an example, Eric Bieschke, head of playlist engineering at Pandora confirms his organization has at least 20 billion thumb ratings from subscribers of the web-based music service. Every twenty four hours, the music company compiles the new recorded actions into the historical database. Actions include thumbs up or down, skipped songs, and new stations built based off of results. This information then undergoes analysis using data mining and integrated filtering tools, to ensure it makes even smarter suggestions for its users going forward. (Mone. 2013) It is this type of machine learning to actively apply learned information within a twenty-four hour period that interests the team at PEGGY.

The value of capturing the data to perform high volume big data analytics is to ensure visitors to PEGGY are doing what the business expects. With session data, translated into a format optimized for near realtime analytics, a system can be built that allows for personalization of PEGGY's website. This personalization will interact with users on a per user basis, greatly enhancing their experience on the site. The main goal of this personalization would be to increase engagement, conversions and average order value (AOV). The AOV is a valuable calculation that represents the sites total revenue divided by the number of orders taken. Analyzing what lead sessions to a purchase greatly impacts predictions and order trends. Uncovering complex dimensions hidden within massive data sets by analyzing pageviews, time on page, unique visits, returning visits, bounce rate, and visitor information is critical to the project. (Marek. 2011)

Current Environment

This evaluation will be comparing to the costs and resources required to upgrade an on premise data warehouse utilizing an Oracle Real Applications Clusters (RAC) System. This is a traditional enterprise relational database system setup that utilizes a traditional extract, transform and load (ETL) process to normalize all data to 3rd normal form. The On-Premise environment presents a limited defined set of reports that can be run on the database. New reports require a global update across the entire database schema when the data does not exist previously. This system does not record the event stream in it's raw form and prevents new reports from capturing historical data. It takes months of resource time to implement a new reporting feature which results in as much as a year long gap in data missing. It also requires lots of testing to ensure the new columns or tables do not have any negative impacts on the existing reports, but this is also dependent on how the new report is implemented.

Finally, upgrading the Oracle Database is extremely costly. The total cost of ownership for the current Oracle implementation is expected to reach \$6,835,200.00 over the next three years. In addition to hardware maintenance and license fees, this system also requires significant employee investments to have experts in house.

Legacy performance (Oracle/ On-Prem)

All of the data is collated utilizing auto_increment primary keys with very few secondary indexes. This has lead to a performance degradation over time as the database tables grow excessively large causing full table scans to take quite some time. All of the data is standardized in third-normal form which is excellent for transactional databases but over time is causing a degradation in the reporting capabilities. Routine nightly processes that ingest all new data have

revealed the dataset is getting so large that indications of significant replication lag post sessionization.

SQL Compliance

Oracle is completely compliant with ANSI SQL and also extends it with PL/SQL. The current system has acquired all of the necessary business intelligence tools to interact with Oracle over ODBC connections. This means that all of the tools use standard SQL, with very little PL/SQL. PL/SQL is used by the database administrator for very specific analytics queries because PL/SQL is the resources speciality. All of these queries could easily be rewritten in standard SQL with little effort, or cost.

Scalability/Performance

Currently the ability to scale and perform in a manner that will meet the business needs is for only another 500 TB of data. At the current rate of data growth, it is estimated the point of no return will pass within six months. There are options within the existing setup to help mitigate this issues such as breaking out the monolithic database into several database instances. There are also alternative indexing techniques and materialized views for the most costly queries that run on the system.

Integrations

Currently the only external integration is a daily batch job of part of the twitter stream that the business feels might be relevant. The process is to normalize the data and insert it into special tables for use by the Marketing Analytics team. This has allowed for the development and purchasing of a set of business intelligence tools to perform analysis on the data. These tools

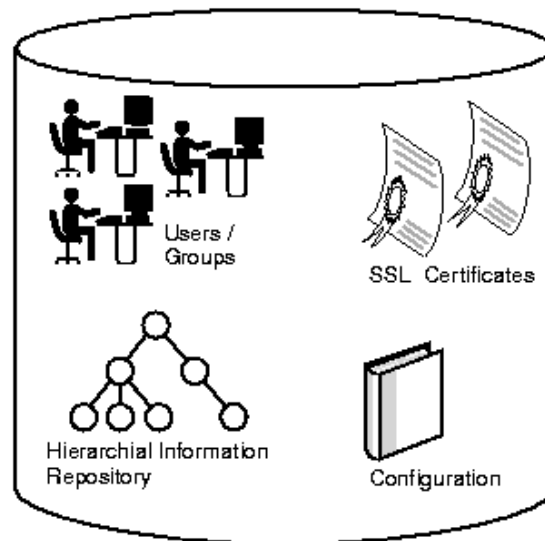
are used by data analysts for every department of the company, with most currently being the marketing department's testing and personalization division. The business intelligence tools in place performs additional testing and recording through the use of third-party applications. The third-party applications validate the internal findings and also uncover key features in the data not provided out-of-the-box. Oracle also allows integration of the data into any dashboard or toolset desired because it can interface through Java Database Connectivity (JDBC) connectors and which allows for the development of new tools around it.

Architecture

The architecture is simple. There are 30 Oracle servers working in an Oracle Real Applications Cluster (RAC) system. This is a centralized database management system maintained by four on staff Database Administrators (DBA). This makes the current Oracle system a single point of failure for all systems. The Oracle servers and RAC licenses are also so expensive it is unaffordable to set up master to slave replication. All backups are done to tape which, as the dataset continues to grow exponentially, requires an exorbitant amount of time and resources to manage. Also, since there is only one set of RAC servers, restoration from tape is rarely tested. This presents the potential for massive data loss at any time. One of the goals of the possible transition to the cloud would be to eliminate the concern for backups and restoration since many of these services have redundancy and recovery built into the systems. In the end, this is a very common data architecture for a small or low growth company. Although the intention of the evaluation is to replace the system as it approaches end of life, it has been a cost effective investment. This project comes at the wake of the next phase of PEGGY's growth and must handle the pre existing and estimated projected needs.

Security

The Oracle Real Application Cluster (RAC) system currently utilizes a role based security backed with authentication and authorization provided by The Lightweight Directory Access Protocol (LDAP) through Microsoft's Active Directory.



source: "Oracle Directory Services (LDAP)," 2000

As displayed in the image above, utilizing LDAP to create specific roles for every employee within the organization controls authentication within a network. This integration is extremely easy to maintain and provides users with granular permissions on databases and tables.

Cost Structure and Resources Required

Currently, based off of the current infrastructure assumptions, an upgrade to a new Oracle database would require approximately 78 (seventy eight) months of development. With three dedicated resources, it is expected to take approximately one year and three months to complete.

Annually, three full time developers along with approximately one and a half full time database administrators would be required.

	Oracle			
Upfront Resource Work	Assume baseline of 72 man-months of application development of application development (Developer salary of \$120,000.00 per year)	\$720,000.00		
	Assume baseline of 6 months of admin effort (Fully-loaded DBA salary of \$120,000.00 per year)	\$60,000.00	\$780,000.00	<-- Year 1
Ongoing Resource Work	Assume baseline of 36 man-months of application development for fully-loaded developer of \$120,000 per yr	\$360,000.00		
	Assumes 1.5 full time DBA's with salaries of \$120,000 per year	\$180,000.00	\$540,000.00	<-- Year 2

source: Appendix, INL 880: Capstone Product Worksheet

The upfront cost for the software licenses and hardware is approximately \$5 million dollars (\$4,580,000.00) with an annual recurring cost (maintenance and support) of \$992,600.00.

Configuration Description			Software: Oracle Database Subscriber Edition & Oracle Real Application Cluster (RAC) Server Hardware: 30 Servers (8 cores/server) w/ 32 GB RAM Storage Hardware: 30 TB SAN (usable)
Upfront Costs	Software Licenses	\$4,230,000.00	\$70,500/RAC core (\$47,500 for Oracle DB Enterprise Edition + \$23,000 for Oracle RAC), 0.5 Xeon Core License Factor, 50% discount off list price.
	Server Hardware	\$120,000.00	8-core servers with 32 GB Ram (\$4,000/server). 30 servers
	Storage Hardware	\$500,000.00	30 TB SAN (usable)
	Total Upfront Costs	\$4,850,000.00	
Annual Ongoing Costs	Software Maintenance and Support	\$930,600.00	22% of license fees
	Server Maintenance and Support	\$12,000.00	10% of hardware purchase price
	Storage Maintenance and Support	\$50,000.00	10% of hardware purchase price
	Total Ongoing Costs	\$992,600.00	

source: Appendix, INL 880: Capstone Product Worksheet

Combining the upfront resource, hardware, and software requirements, it is estimated that over the course of three years, the total cost of ownership will be \$8,695,200.00.

	Oracle Database Subscriber Edition & Oracle Real Application Cluster (RAC) Server Hardware: 30 Servers (8 cores/server) w/ 32 GB RAM Storage Hardware: 30 TB SAN (usable)
Total Upfront Resource Costs	\$780,000.00
Total Ongoing Resource Costs	\$540,000.00
Total Upfront Infrastructure Cost	\$4,850,000.00
Total Ongoing Infrastructure Cost	\$992,600.00
Total Year 1	\$5,630,000.00
Total Year 2	\$1,532,600.00
Total Year 3	\$1,532,600.00
3 Year TCO	\$8,695,200.00

source: Appendix, INL 880: Capstone Product Worksheet

Advantages of Legacy infrastructure

Normalization of all of the data allows all engineers to easily understand the data model and write reports against the data. Enforcing a strict structure also allows the team to manage multiple requests simultaneously. This legacy architecture also opens up a large talent pool when compared to a newer cloud based big data offerings. Proper staffing is extremely important for development of new reports and internal applications that utilize the data. It allows any type of software engineer with relational database experience to develop new tools. However, the limits it imposes on scalability may not make this a great tradeoff.

Disadvantages of Legacy infrastructure

The need to standardize reports in order to reduce the cost of implementation is holding the business back. The backup and failure scenarios for the existing infrastructure are serious points of concern. There is only one single point of failure which means that, if the Oracle system goes down, there would be a business significant outage for all internal reporting applications. This would have the impact of delaying site advancements, marketing campaigns and thus would have a direct impact on revenue. If the Oracle RAC were to have a critical hardware failure now, there is no guarantee for the safety of the data. There is limited protection against single hard drive failure but not against the outage of an entire set of hard drives. If, for example, the air conditioning stopped working in the server room and it could not be repaired in time, major damage could be done to the physical servers. There is no guarantee that the restore from tape would be 100% effective and data loss from the time after the previous backup to current would be inevitable. This is not a position any IT department, let alone a quickly growing one, wants to be in it.

Evaluation of Options**On-Premise vs The Cloud**

On-Premise and Cloud Computing are developed with very different frameworks. On-Premise is most commonly associated with static models that are incapable of change, whereas cloud computing is widely praised for its non-linear dynamic models capable to scale up or down as needed. Cloud computing means accessing data, applications, storage, and computing power over the web rather than on the hard drives of premise based machines. (Watson.2014)

Cloud computing, as it relates to infrastructure, enables systems that are themselves adaptive and dynamic to handle the increase (or decrease) in demand and automatically optimize while utilizing the extensive resources available. Vendors offering an Infrastructure as a Service model, like Amazon, maintain computer servers, storage servers, communication infrastructure, and all common data center services. A data center is a large facility where the hardware, uninterrupted power supply, access control, and communication facility are located. It is at these data centers where the hosted systems and application software rests. Additionally, IaaS solutions, in most cases, offer multi tenanted, which means the cloud vendors offer a public cloud solution where a single instance is shared to multiple users. Based on leaders in IaaS offerings, Amazon has been in the business the longest and first started with Elastic Compute Cloud (EC2). (Rajaraman. 2014) To clarify, EC2 is an interface that delivers a web-based server environment and gives users full control to provision any number of servers in minutes regardless of scale or capacity. Other larger players in the IaaS provider space are Rackspace, IBM (SmartCloud+), Microsoft, and Google. All these providers offer various types of virtualized systems to scale to the programming needs.

Today, cloud computing in the enterprise space is widely known for the adoption of the on-demand, pay-as-you-go service rather than the traditional on-premise locally stored, managed, and operated model. There is a vast array these types of a service offerings such as software as a service (SaaS), platform as a service (PaaS), and desktops as a service (Daas). This report will be focused primarily on the cloud offerings of infrastructure as a service (IaaS).

Cloud investments as a whole have grown 19% over 2012 and, in the next 1 to 3 years, 35% of business/ data analytics projects will go to the cloud. (IDG_Enterprise. 2014) Further, 24% of IT budgets slated for 2015 are devoted to cloud solutions, 28% of this is for IaaS and

18% for PaaS. (Columbus. 2014) Cloud solutions are rapidly improving time-to-market capabilities while also reducing the total cost of ownership.

It has been recently reported by Gartner that of all the IaaS offerings, Amazon Web Services far outpaces the competition of computing power when compared to Microsoft, Google, and IBM.

Figure 1 - Gartner 2014 Magic Quadrant for Cloud Infrastructure as a Service.



source: "Gartner's Magic Quadrant," 2014

This magic quadrant evaluated current cloud based IaaS in the context of hosting a data center in the cloud. These types of IaaS solutions allow for the user to still retain most IT control

such as governance and security and the ability to run both new and legacy workloads. (Gartner. 2014)

As an example, a one thousand terabyte dataset required the PEGGY team to perform a set of four operations on every piece of data within the data set. On the system currently being used at Peggy (read: on-premise big iron solution) a massive server and storage system would be necessary along with a fibre connection in order to fully maximize bandwidth. The task certainly can and will be completed but the computing pace would likely be a deterrent. This is known as I/O bound, because the time it takes to complete a computation is determined by the period spent waiting for input/output operations to be completed. (Turkington. 2013) Due to the size and complexity of the datasets, more time is spent requesting the data than processing it. Consider the Pandora example.

Alternatively, cloud based solutions remove the tasks relevant to infrastructure, and instead, focus on either utilizing pre-built (ie: public cloud vendors like Amazon Web Services) or assigning developers to build cloud-based applications (ie: open source) to perform the same task. Both open source and IaaS systems handle the cluster mechanics transparently. These models (ie: open source and IaaS) allows the developers or data analysts to think in terms of the business problem. (Turkington. 2013) Further, Google's parallel cloud-based query service Dremel has the capability to "scan 35 billion rows without an index in tens of seconds. (Sato.2012) Dremel is capable of doing this by parallelizing each query and running it on tens of thousands of servers simultaneously. This type of technology eliminates ongoing concerns regarding processing speed in proportion to CPU speed (ie: I/O bound) entirely. As pointed out by Google, no two clouds are the same and they are offered as both bundled and a la carte purchasing options. (Ward. January 28, 2015) Pay as you go services like IaaS require a different

7/20/2015 INL-880 - Capstone Proposal: McGinley & Etter -Final Draft

mindset. Rather than the upfront capital expenditure of massive ironbound infrastructure, cloud system offer a pay-as-you-go model, and transition computing power, such as storage and analytics into an operational expenditures. Cloud-based vendors, such as Google and Amazon, inherit the responsibility for system health and support. Further, additional storage and hardware costs are no longer a consideration. (Hertzfeld. 2015)

These monthly recurring operational costs require a new frame of mind in order to budget. For the purpose of determining an adequate cloud service to replace the current Oracle database, usage hours per day/ month/ year; instance cost; number of servers; operating system (o/s), central processing units (CPUs) often referred to as number of cores; random access memory (RAM); solid state drive (SSD) or hard disk drive (HDD); regions/ zones/ collocations; upfront costs in addition to monthly recurring fees; reserved (ie: annual or multiyear) vs on demand commitment/ agreement terms all need to be considered. In addition to cost effectiveness, and separate from development, programming, and administration, the cloud services remove the tasks of deploying, managing and upgrading infrastructure to scale.

Open Source vs Infrastructure as a Service

There is a growing argument between cloud services regarding whether or not to favor open standards due to the diversity and capability. Open source software is always available at no cost which is reason that quality, in many stages, is uncomparable to the turnkey solutions provided by proprietary services such as IaaS. (Leoncini. 2011)

Amazon and Google offer both open source and private cloud offerings. These tools are helping organizations essentially rent computers, apps and storage in remote data centers via the web to build their own private, internal cloud. (Krause. 2002) Similarly, both Google and

Amazon deliver a web-portal for users to rent servers for as little or long as needed in a utility-like model.

Amazon and Google collectively started the wave of low-cost broadband communications offerings with unprecedented speed and storage capacities of computers with on-demand costs.. Both organization quickly became the two leading competitors of cloud services between 2004 and 2006.

The computing facility Amazon was using for it's online book and shopping store was not operating at full utilization. (less than 10%). This was seen as a business opportunity to sell the excess computing infrastructure. In 2006 Amazon started Amazon Web Services which sold computing infrastructure on demand using the Internet for communication. (Rajaraman. 2014)

Similarly, Google was the leader as a free search engine and required a large computing infrastructure to cater to the most optimal search speed expected. In 2004, Google released a free email service, GMail, for all its customers using this infrastructure and in 2006 expanded its offerings to include free office productivity suite called Google Docs with 2GB free disk space. Similar to Amazon, Google recognized a business opportunity to sell excess hardware capacity and started Google compute engine as a paid cloud service in 2012. (Rajaraman. 2014)

While Google's search engine was evolving, the team at Google needed to implement hundreds of specific computations in order to process large amounts of raw data (crawled documents, web request logs, etc.). In order to handle the increasing demand of the growing user base, they needed to determine a way to "parallelize the computation, distribute the data, and handle failures conspire to obscure the original simple computation with large amounts of complex code to deal with these issues." (Dean. 2004) Once Google discovered a solution to their problem, they released two academic papers which described the platform to process data

highly efficiently on very large scale. The papers discussed two technologies, Google File System (GFS) and MapReduce.

MapReduce is a programming model that was created to deliver an interface that enables automatic parallelization and distribution of large-scale computations and high performance on large clusters of commodity PCs. (Dean. 2004) Google File System is a technology that distributes massive amounts of data across thousands of inexpensive computers. This technology allows Google to support large-scale data processing workloads to commodity, or rather, traditional hardware. Further, the system is fault tolerant through constant monitoring, replicating crucial data, and fast automatic recovery. (Ghemawat. 2003) Google expects that all machines will fail, so building failure into their model allowed for them to dramatically reduce the infrastructure cost while achieving high capacity computing.

These two papers resulted in the creation of several open source software offerings, most notably, Apache Hadoop which also has two offerings. The Hadoop Distributed File System (HDFS) shares and stores enormous datasets among thousands of inexpensive pieces of hardware. Hadoop MapReduce takes the information from HDFS and computes the separated dataset on independent machines and processing power. The two combined offer a compelling storage and processing offering in the cloud.

Since the release of the originating documents, data storage systems available for reporting and analytics has grown exponentially. Systems like Amazon Redshift offer data warehousing in a traditional data center format which allows for manual configuration and administration without the need to purchase and maintain hardware. However, data warehousing in the cloud services such as Google's BigQuery offers an untraditional service by offering

elastic storage, network, and computing capabilities without any additional provisioning or administration through automatic scaling. (Ward. June #, 2015)

Standard data types captured for sessionized data

There is some standard information organizations want to capture about users, both anonymous and known. When a person converts from anonymous to known, organizations start an event so that they can match the user's anonymous history with their known history. Some of the basic data we want to capture is What products (ie: product id) and variations of the product (sku) a user has seen, perhaps even what images they hovered or lingered on. They want to know when users add and remove things to their cart, what they actually buy. They want to capture the User Agent (UA) string from the browser so that it can be can determine what platforms the user has and engages the site from. They will also want to track IPs and do per request geoip lookups and record the result so we know where the user was accessing the site from. All of this information will allow them to run the normal ecommerce analytics queries and understand more about customers. It allows them to segment the population of the site into groups they know and understand and calculate their customer lifetime value (CLV), which helps us understand where to put the marketing efforts for the biggest positive impact on the company.

Standard analytic queries used in e-commerce

All businesses utilize Key Performance Indicators (KPI), which are measurable values that demonstrate how effectively an organization is at achieving it's objectives. (Rouse. 2006) A large majority of ecommerce companies care about the same types of analytics queries, which are the KPIs for these organizations. This is true of PEGGY also. The primary indicators of concern are Average Order Value (AOV), conversion rate, the average number of pageviews,

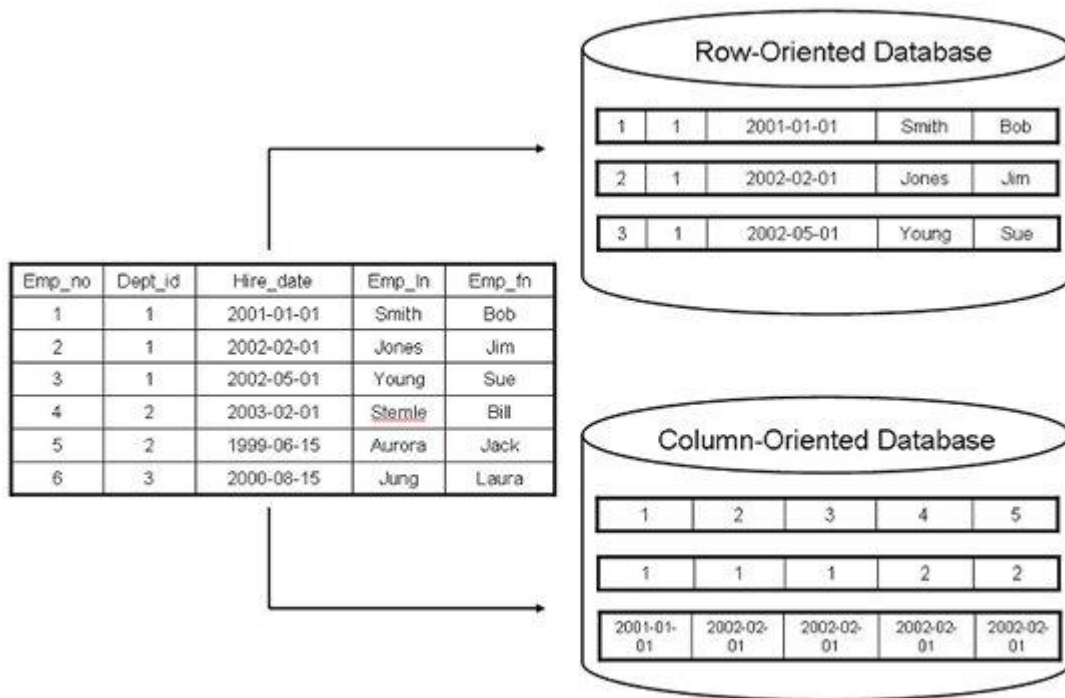
and the number of abandoned carts. These KPIs help the marketing team determine the top level input to the organization.

Using this, the marketing and inventory teams look at what products and product variations users are engaging with the most, to determine reorder information and to give them ideas for new products. Marketing and the IT department are also curious about platform information to determine where bugs or issues with the user interface may be interfering with the user engagement. The marketing department leverages analytical reporting on demographic information such as user locations, site traffic, bounce rate, lift on targeted campaigns, and a lot of other queries to determine what specific efforts are taking have a positive effect on who. The more granular and detailed the reports are, the more obvious the impact of small changes are on KPIs for the company.

Analysis of Products

Amazon RedShift

Amazon Redshift is a Columnar Database designed for Petabyte scale provided as a hosted service. A column database stores the data contained in it to disk in a different manner from traditional databases.



source: Moore. (2011)

David Raab in his article “How to Judge a Columnar Database” has an excellent description of how they differ from traditional databases, “As the name implies, columnar databases are organized by column rather than row: that is, all instances of a single data element

(say, Customer Name) are stored together so they can be accessed as a unit. This makes them particularly efficient at analytical queries, such as list selections, which often read a few data elements but need to see all instances of these elements. In contrast, a conventional relational database stores data by rows, so all information for a particular record (row) is immediately accessible.” (Raab. 2007)

Amazon Redshift converts the data to columnar storage automatically and in the background. Amazon has determined this methodology will increase storage efficiency substantially for tables that have large numbers of columns and very large row counts. Additionally, Amazon notes that since each block contains the same type of data, they can apply a compression scheme specific to the column data type, and reduce disk space and I/O further.

This impacts memory as well as, due to the need to only pull data within specific rows or columns, memory is saved by selecting the individual blocks as opposed to the entire row or column. When compared to typical OLTP or relative data warehouse query, Redshift is capable of utilizing a fraction of the memory process information.. (“Database Developer Guide”,2015) Redshift also utilizes the capabilities of a hosted service to increase query performance. When a Redshift cluster is initiated, the administrator is allocated special servers within the AWS infrastructure. A notable feature is Redshift offers solid state drives (SDD) rather than standard hard drives (HDD). The instances allocated also utilize high performance memory hardware, which allows them to store large amounts of data in memory and quickly fetch it from disk. Combined together the specialized hardware and software allows Amazon Redshift to store Petabytes of data and quickly run analytical queries on it.

SQL Compliance

Amazon Redshift has significant ANSI SQL compliance. Amazon in fact states “Many of your queries will, work with little or no alterations from a syntax perspective.” There are really only a small number of functions that Redshift does not support including “convert()” and “substr()” and generally these are not supported for performance reasons. Redshift also adds some functions to help optimize the performance of queries on extremely large datasets. In fact all of the additions and constraints added to the SQL compliance of Redshift are around the performance on large datasets. For example if we look back at convert and substr, these are removed because they would have to be executed on every row of a table being queried, which is highly non performant at petabyte scale. The other main difference between standard SQL and Redshift is the idea of distribution keys and sort keys. These keys tell Redshift how to optimally split data across it’s hard drives and nodes for future querying. Primary keys and foreign keys can be defined in Redshift but it expects that the referential integrity to be enforced by the program inserting data, and the database itself will allow duplicates, and bad references. Again the reason that Redshift does not enforce these keys by default is for performance because large table scans would have to occur in some cases to enforce these keys, destroying insert performance. In fact, Amazon suggests never doing single row INSERTs into Redshift. The preferred method is to use bulk inserts from Amazon’s Simple Storage Service (S3) or a file located on a server. This is because individual inserts often cause more work for the server during distribution and sorting as opposed to bulk inserts which can be optimized to insert.

Multi-row inserts improve performance by batching up a series of inserts. The following example inserts three rows into a four-column table using a single INSERT statement.

This is still a small insert, shown simply to illustrate the syntax of a multi-row insert.

source: “use a multi-row insert”, 2015

Amazon’s recommendation to only use batch inserts is a prime example why Redshift should not be used as a transactional database but instead exclusively as a data warehouse for analytics. One other final note of some importance is that command line connections to Redshift occur with an older version of the PostgreSQL command line tool. This let’s us know that Redshift has a programmatic basis in PostgreSQL of some type. This is important because it also gives us an idea about what kinds of drivers will work with Redshift for programmatic access.

Performance and Scalability

Amazon Redshift is designed to be highly performant for queries on datasets up to petabytes in size. Amazon supports petabyte datasets with a Redshift cluster, but there are limits placed on the max size of a cluster you can have based on what type of cluster you setup in Amazon. There are four types of nodes that a Redshift cluster can have, Amazon provides the following tables for basic node type information.

Dense Storage Node Types

Node Size	Node Limits	Storage Capacity per Node	Maximum Storage Capacity per Cluster
dw1.xlarge	1 to 32	2 TB hard disk drive (HDD) storage	64 TB
dw1.8xlarge	2 to 128	16 TB hard disk drive (HDD) storage	2 PB

source: awsdocumentation. 2015 About Clusters and Nodes

Dense Compute Node Types

Node Size	Node Limits	Storage Capacity per Node	Maximum Storage Capacity per Cluster
dw2.large	1 to 32	160 GB solid state drive (SSD) storage	5.12 TB
dw2.8xlarge	2 to 128	2.56 TB solid state drive (SSD) storage	326 TB

source: awsdocumentation. 2015 About Clusters and Nodes

These node types put the max size of a cluster, utilizing the node size entitled dw1.8xlarge, as noted in the chart above, at 256 Petabytes. This well exceeds the requirements for storage for the long term. When you have a cluster of any size, Amazon uses the distribution keys to distribute data across the cluster of nodes you have set up. It is important to choose a distribution key that will help Amazon easily spread all of your data evenly across your cluster, because then each node can work effectively at filtering data in response to queries. More complex queries, for example those with a 'join' or a 'group by' will require data to be moved around the cluster and the distribution of data can help make sure that smaller amounts of data are transferred to the leader node for locality. The leader node is a free service that Amazon provides that "receives queries from client applications, parses the queries and develops execution plans, which are an ordered set of steps to process these queries." (Amazon Web Services. "Redshift FAQ's.")

Many optimizations also occur when a user sends a SQL query to Redshift. Specifically since the data storage format is specific and custom, a key part of the query engine can be written efficiently. Specifically the SQL query optimizer analyzes the statement and Redshift then creates a small C++ executable that is distributed to all the nodes. Since the storage format of Redshift is so very specific and explicit the application is then executed on all the nodes and the data is pulled from storage on that node and then decisions are made about what to do with it.

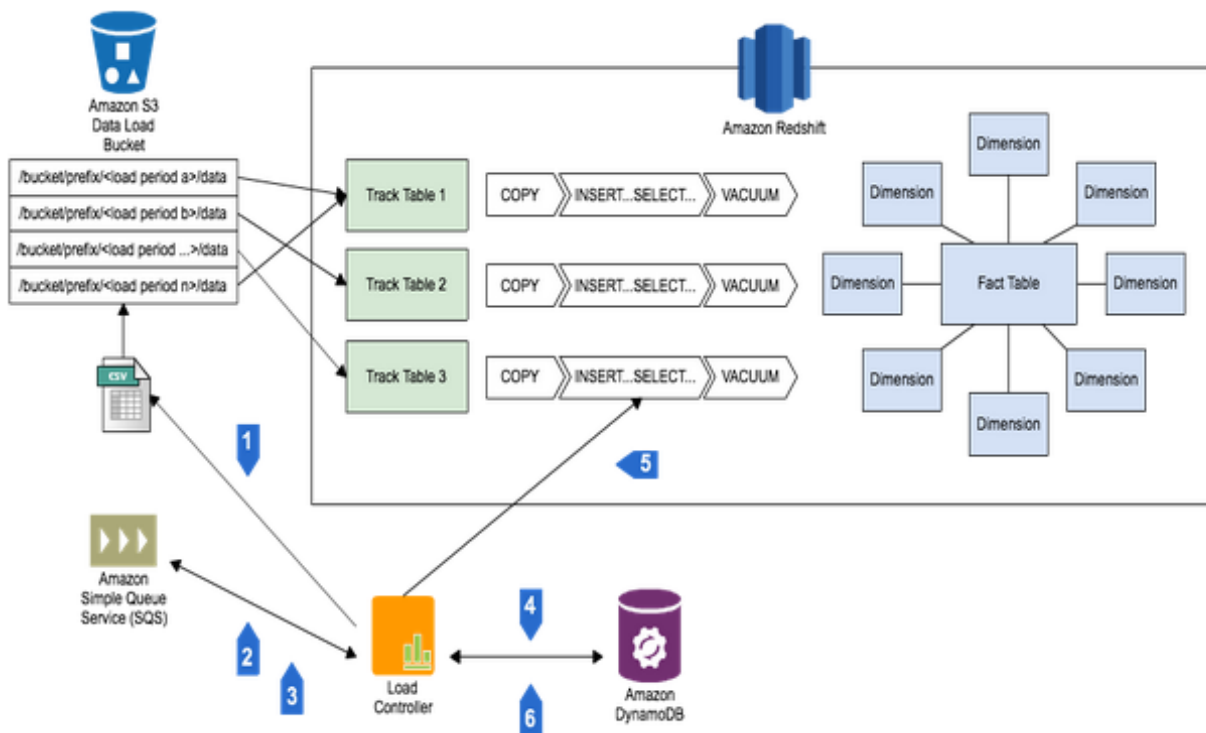
Some things that can happen with this data include, sending it all to the leader for further filtering. This is in fact a performance bottleneck, which Redshift will explain in query analysis by providing you with the DS_BCAST_INNER keyword that provides the administrator a copy of the entire inner table which is broadcasted to all the compute nodes. (“analyzing the query plan,” 2015.) Amazon also include queries like DS_BCAST_INNER, which tells you that all data is going to one node for joining and querying, which is extremely network and memory intensive. Other hits include DS_DIST_ALL_INNER which “Indicates that all of the workload is on a single slice.” and DS_DIST_BOTH which “Indicates heavy redistribution.” Redshift also provides tables that log both queries waiting to be run and those that have recently been run so that users can do analytics on how long their queries are taking and then look for performance gains in these queries. In fact Redshift provides several analysis tools for users to find bottlenecks in their queries. Overall, Redshift provides us with the tools and capabilities to maintain performance and to scale the data set easily into the Petabyte range. As for speed, Stefan Bauer, author of Getting Started with Amazon Redshift noted, "We took Amazon Redshift for a test run the moment it was released. It's fast. It's easy. Did I mention it's ridiculously fast? We've been waiting for a suitable data warehouse at big data scale, and ladies and gentlemen it's here. We'll be using it immediately to provide our analysts an alternative to Hadoop. I doubt any of them will want to go back." (Bauer. 2013)

Additionally, Amazon explains that Redshift allows segmentation of workload. Batch operations and reporting like data exploration can be separated from less resource-intensive queries. In turn, this type of manual configuration will boost overall performance speed. (Keyser. 2015)

An example of segmentation is as follows:

7/20/2015

INL-880 - Capstone Proposal: McGinley & Etter -Final Draft



source: “optimizing star schemas on Redshift,” 2015

Integrations

Amazon Redshift offers several integrations with multiple data extract, transform, and load (ETL) and business intelligence (BI) reporting, data mining, and analytics tools. Redshift’s design around PostgreSQL which, in effect, enables most SQL client applications to work and function with minimal disruption or change. (“Database Developer Guide,” 2015) . Redshift also includes JDBC and ODBC support which enables common tools such as Tableau and Looker to function with minimal change. The ability to integrate with all these tools and scale to support large data sets makes Amazon’s Redshift product an excellent data store for business analytics teams. Infoworld.com has a quote from the launch of Amazon Redshift showing the importance

of this compability “AWS CTO Werner Vogels blogged that ‘Amazon Redshift enables
7/20/2015

INL-880 - Capstone Proposal: McGinley & Etter -Final Draft

customers to obtain dramatically increased query performance when analyzing datasets ranging in size from hundreds of gigabytes to a petabyte or more, using the same SQL-based business intelligence tools they use today.”(Lampitt. 2012) Utilizing Amazon Redshift would enable the PEGGY system to keep all of the investments in Analytics Visualization and Business Intelligence tools for years to come. Redshift will also allow these tools to remain relevant for a much longer time, by scaling the data to a size the Oracle RAC would be incapable of handling.

Architecture

Amazon states their solution offers ten time the performance capabilities of traditional on-premise data warehousing and analytics solutions:

The biggest recent Big Data announcement in that field, SAP’s HANA, an in-memory high power database management platform that app developers are rushing to design to, now seems eclipsed by Redshift. The irony is that SAP is touting HANA because it offers a powerful solution at budget price because it can run on the Amazon cloud: ‘just’ \$300,000. That’s impressive performance for the price – but now Redshift can give you most of that for one third of one percent of SAP’s price. (Peters. 2013)

In addition to utilizing columnar data storage, Redshift achieves efficient storage and optimum query performance through a combination of massively parallel processing and very efficient, targeted data compression encoding schemes.

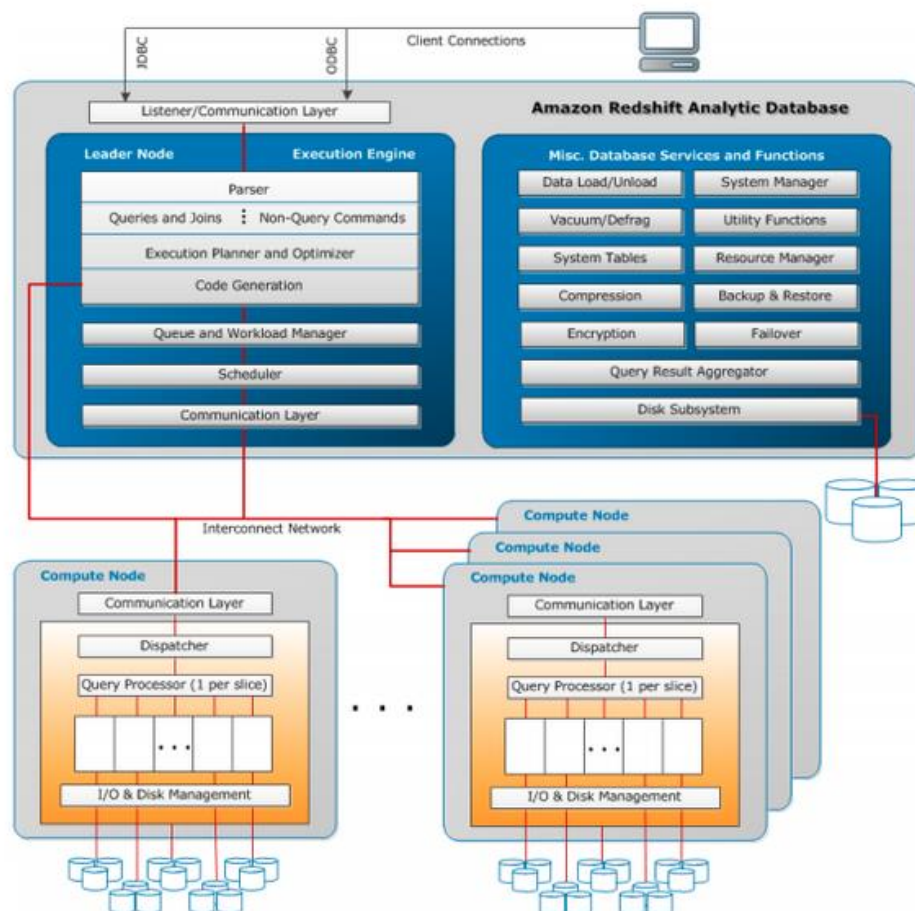
According to Peter Scott, of Rittman Mead Consulting:

A key point of difference between Amazon Redshift and Oracle is in how the data is stored or structured in the database. An understanding of this is vital in how to design a performance data warehouse. With Oracle we have shared storage (SAN or local disk)

attached to a pool of processors (single machine or a cluster); however, Redshift uses a share-nothing architecture, that is the storage is tied to the individual processor cores of the nodes. As with Oracle, data is stored in blocks, however the Redshift block size is much larger (1MB) than the usual Oracle block sizes; the real difference is how tables are stored in the database, Redshift stores each column separately and optionally allows one of many forms of data compression. Tables are also distributed across the node slices so that each CPU core has its own section of the table to process. In addition, data in the table can be sorted on a sort column which can lead to further performance benefits.

(Scott. 2014)

As noted, Amazon Redshift is a relational database management system (RDBMS) and is compatible with most common on premise applications. Although it provides similar functions such as inserting and deleting data, Amazon Redshift is optimized to quickly scale up or down in order to deliver high-performance analysis and reporting of very large datasets.



source: "Database Developer Guide", 2015

As indicated by the image above, the Redshift primary infrastructure is centered around clusters, which represent a collection of one or more compute nodes. Each cluster could contain one or multiple databases. When provisioning a multiple compute node cluster, an additional leader node is created to communicate between external client communications and the compute nodes. The leader node will communicate exclusively with the on premise SQL client. The queryable data is then split across all compute nodes on the cluster in an Amazon specific manner to optimize query performance. The compute nodes each have their own dedicated CPU, memory, and attached disk storage which is predetermined based on the node type. However,

increasing the compute and storage capacity of a cluster by increasing the number of nodes or upgrading the node type can be done at any time. (“Database Developer Guide”, 2015)

Disaster recovery is also maintained by Amazon. “Amazon Redshift replicates all your data within your data warehouse cluster when it is loaded and also continuously backs up your data to S3. Amazon Redshift always attempts to maintain at least three copies of your data (the original and replica on the compute nodes and a backup in Amazon S3). Redshift can also asynchronously replicate your snapshots to S3 in another region for disaster recovery.” (Amazon Redshift FAQs). This is extremely important as it prevents a team from having to exert any effort to guarantee data safety, and allows extremely quick recovery from a failure.

Security

AWS has in the past successfully completed multiple SAS70 Type II audits, and now publishes a Service Organization Controls 1 (SOC 1), Type 2 report, published under both the SSAE 16 and the ISAE 3402 professional standards as well as a Service Organization Controls 2 (SOC 2) report. In addition, AWS has achieved ISO 27001 certification, and has been successfully validated as a Level 1 service provider under the Payment Card Industry (PCI) Data Security Standard (DSS). In the realm of public sector certifications, AWS has received authorization from the U.S. General Services Administration to operate at the FISMA Moderate level, and is also the platform for applications with Authorities to Operate (ATOs) under the Defense Information Assurance Certification and Accreditation Program (DIACAP). (“AWS Cloud Security,” 2015) Amazon has undergone numerous additional compliance audits in order to assure their customers the cloud infrastructure meets the needs surrounding security and protection.

Here is a list of all of the relevant security audits and programs Amazon has undergone that are relevant to E Commerce and organizations headquartered in the United States:

Audit/ Program	Explanation
PCI DSS Level 1	AWS is Level 1 compliant under the Payment Card Industry (PCI) Data Security Standard (DSS). Customers can run applications on their PCI-compliant technology infrastructure for storing, processing, and transmitting credit card information in the cloud.
FedRAMP (SM)	AWS has achieved two Agency Authority to Operate (ATO) under the Federal Risk and Authorization Management Program (FedRAMP) at the Moderate impact level. FedRAMP is a government-wide program that provides a standardized approach to security assessment, authorization, and continuous monitoring for cloud products and services up to the Moderate level.
HIPPA	AWS enables covered entities and their business associates subject to the U.S. Health Insurance Portability and Accountability Act (HIPAA) to leverage the secure AWS environment to process, maintain, and store protected health information. Additionally, AWS, as of July 2013, is able to sign business associate agreements (BAA) with such customers.
SOC 1/ ISAE 3402	Amazon Web Services publishes a Service Organization Controls 1 (SOC 1), Type II report. The audit for this report is conducted in accordance with AICPA: AT 801 (formerly SSAE 16) and the International Standards for Assurance Engagements No. 3402 (ISAE 3402). This audit is the replacement of the Statement on Auditing Standards No. 70 (SAS 70) Type II report. This dual-standard report can meet a broad range of auditing requirements for U.S. and international auditing bodies.
DIACAP and FISMA	AWS enables US government agencies to achieve and sustain compliance with the Federal Information Security Management Act (FISMA). The AWS infrastructure has been evaluated by independent assessors for a variety of government systems as part of their system owner's approval process. Numerous Federal Civilian and Department of Defense (DoD) organizations have successfully achieved security authorizations for systems hosted on AWS in accordance with the Risk Management Framework (RMF) process defined in NIST 800-37 and DoD Information Assurance Certification and Accreditation Process (DIACAP).
Dod CSM Levels 1-2, 3-5	The Department of Defense (DoD) Cloud Security Model (CSM) provides a formalized assessment and authorization process for cloud service providers (CSPs) to gain a DoD Provisional Authorization, which can subsequently be leveraged by DoD customers. A Provisional Authorization under the CSM provides a reusable certification that attests to our compliance with DoD standards, reducing the time necessary for a DoD mission owner to assess and authorize one of their systems for operation on AWS.
SOC 2	In addition to the SOC 1 report, AWS publishes a Service Organization Controls 2 (SOC 2), Type II report. Similar to the SOC 1 in the evaluation of controls, the SOC 2 report is an attestation report that expands the evaluation of controls to the criteria set forth by the American Institute of Certified Public Accountants (AICPA) Trust Services Principles. These principles define leading practice controls relevant to security, availability, processing integrity, confidentiality, and privacy applicable to service organizations such as AWS.
SOC 3	AWS publishes a Service Organization Controls 3 (SOC 3) report. The SOC 3 report is a publicly-available summary of the AWS SOC 2 report. The report includes the external auditor's opinion of the operation of controls (based on the AICPA's Security Trust Principles included in the SOC 2 report), the assertion from AWS management regarding the effectiveness of controls, and an overview of AWS Infrastructure and Services.
ISO 27001	AWS is ISO 27001 certified under the International Organization for Standardization (ISO) 27001 standard. ISO 27001 is a widely-adopted global security standard that outlines the requirements for information security management systems. It provides a systematic approach to managing company and customer information that's based on periodic risk assessments. In order to achieve the certification, a company must show it has a systematic and ongoing approach to managing information security risks that affect the confidentiality, integrity, and availability of company and customer information.
ISO 9001	ISO 9001:2008 is a global standard for managing the quality of products and services. The 9001 standard outlines a quality management system based on eight principles defined by the International Organization for Standardization (ISO) Technical Committee for Quality Management and Quality Assurance.

	<p>They include:</p> <ul style="list-style-type: none"> ● Customer focus ● Leadership ● Involvement of people ● Process approach ● System approach to management ● Continual Improvement ● Factual approach to decision-making ● Mutually beneficial supplier relationships
MPAA	The Motion Picture Association of America (MPAA) has established a set of best practices for securely storing, processing, and delivering protected media and content. Media companies use these best practices as a way to assess risk and security of their content and infrastructure. AWS has demonstrated alignment with the MPAA Best Practices and AWS infrastructure is compliant with all applicable MPAA infrastructure controls.
CJIS	In the spirit of a shared responsibility philosophy AWS has created a Criminal Justice Information Services (CJIS) Workbook in a security plan template format aligned to the CJIS Policy Areas. This Workbook is intended to support our partners documenting their alignment to CJIS security requirements.
FIPS 140-2	The Federal Information Processing Standard (FIPS) Publication 140-2 is a US government security standard that specifies the security requirements for cryptographic modules protecting sensitive information. To support customers with FIPS 140-2 requirements, SSL terminations in AWS GovCloud (US) operate using FIPS 140-2 validated hardware.
Section 508/ VPAT	Section 508 was enacted to eliminate barriers in information technology, to make available new opportunities for people with disabilities, and to encourage development of technologies that will help achieve these goals. The law applies to all Federal agencies when they develop, procure, maintain, or use electronic and information technology. Under Section 508 (29 U.S.C. ' 794d), agencies must give disabled employees and members of the public access to information that is comparable to the access available to others. Amazon Web Services offers the Voluntary Product Accessibility Template (VPAT) upon request.
FERPA	The Family Educational Rights and Privacy Act (FERPA) (20 U.S.C. § 1232g; 34 CFR Part 99) is a Federal law that protects the privacy of student education records. The law applies to all schools that receive funds under an applicable program of the U.S. Department of Education. FERPA gives parents certain rights with respect to their children's education records. These rights transfer to the student when he or she reaches the age of 18, or attends a school beyond the high school level. Students to whom the rights have transferred are "eligible students."
ITAR	The AWS GovCloud (US) region supports US International Traffic in Arms Regulations (ITAR) compliance. As a part of managing a comprehensive ITAR compliance program, companies subject to ITAR export regulations must control unintended exports by restricting access to protected data to US Persons and restricting physical location of that data to the US. AWS GovCloud (US) provides an environment physically located in the US and where access by AWS Personnel is limited to US Persons, thereby allowing qualified companies to transmit, process, and store protected articles and data subject to ITAR restrictions.
CSA	In 2011, the Cloud Security Alliance (CSA) launched STAR, an initiative to encourage transparency of security practices within cloud providers. The CSA Security, Trust & Assurance Registry (STAR) is a free, publicly accessible registry that documents the security controls provided by various cloud computing offerings, thereby helping users assess the security of cloud providers they currently use or are considering contracting with. AWS is a CSA STAR registrant and has completed the Cloud Security Alliance (CSA) Consensus Assessments Initiative Questionnaire (CAIQ). This CAIQ published by the CSA provides a way to reference and document what security controls exist in AWS's Infrastructure as a Service offerings. The CAIQ provides a set of over 140 questions a cloud consumer and cloud auditor may wish to ask of a cloud provider.

source: "AWS Compliance," 2015

Amazon Redshift security is maintained by both Amazon Identity and Access

Management (IAM) and users that can be setup in the database, as is common with MySQL and

other databases. Access can also be restricted utilizing Security Groups. These security groups take CIDR blocks to restrict all port access to a server by IP; this is much like you would find when using IP Tables on a standard Linux server. All access to the Redshift servers is also monitored and logged to Amazon cloud watch.

In addition, Amazon Redshift supports Amazon Virtual Private Cloud (Amazon VPC), SSL, AES-256 encryption and Hardware Security Modules (HSMs) to protect data in transit and at rest.

Sign-in credentials	— Access to your Amazon Redshift Management Console is controlled by your AWS account privileges. For more information, see Sign-In Credentials.
Access management	— To control access to specific Amazon Redshift resources, you define AWS Identity and Access Management (IAM) accounts. For more information, see Controlling Access to Amazon Redshift Resources.
Cluster security groups	— To grant other users inbound access to an Amazon Redshift cluster, you define a cluster security group and associate it with a cluster. For more information, see Amazon Redshift Cluster Security Groups.
VPC	— To protect access to your cluster by using a virtual networking environment, you can launch your cluster in a Virtual Private Cloud (VPC). For more information, see Managing Clusters in Virtual Private Cloud (VPC).
Cluster encryption	— To encrypt the data in all your user-created tables, you can enable cluster encryption when you launch the cluster. For more information, see Amazon Redshift Clusters.
SSL connections	— To encrypt the connection between your SQL client and your cluster, you can use secure sockets layer (SSL) encryption. For more information, see Connect to Your Cluster Using SSL.
Load data encryption	— To encrypt your table load data files when you upload them to Amazon S3, you can use either server-side encryption or client-side encryption. When you load from server-side encrypted data, Amazon S3 handles decryption transparently. When you load from client-side encrypted data, the Amazon Redshift COPY command decrypts the data as it loads the table. For more information, see Uploading Encrypted Data to Amazon S3.
Data in transit	— To protect your data in transit within the AWS cloud, Amazon Redshift uses hardware accelerated SSL to communicate with Amazon S3 or Amazon DynamoDB for COPY, UNLOAD, backup, and restore operations.

source: "Amazon Redshift Security Overview," 2015

Cost Structure

The cost structure behind Redshift is relatively simple. When spinning up an instance of Redshift, you can choose between On-Demand or Reserved Instances. Additionally, there's an option to choose between dense storage (DS) or dense compute (DC) nodes. The difference between dense compute and dense storage is, when creating a data warehouse, dense storage nodes is more focused on utilizing hard disk drives for very large datasets and dense compute is

for high capacity for performance utilizing fast CPUs and and RAM through solid-state disks (SSDs).

Region: US East (N. Virginia) ▾						
	vCPU	ECU	Memory (GiB)	Storage	I/O	Price
DW1 - Dense Storage						
dw1.xlarge	2	4.4	15	2TB HDD	0.30GB/s	\$0.850 per Hour
dw1.8xlarge	16	35	120	16TB HDD	2.40GB/s	\$6.800 per Hour
DW2 - Dense Compute						
dw2.large	2	7	15	0.16TB SSD	0.20GB/s	\$0.250 per Hour
dw2.8xlarge	32	104	244	2.56TB SSD	3.70GB/s	\$4.800 per Hour

source: “Amazon Redshift Pricing,” 2015

The pay-as-you-go offering known as on-demand instances let you pay for compute capacity by the hour with no long-term commitments. This frees you from the costs and complexities of planning, purchasing, and maintaining hardware and transforms what are commonly large fixed costs into much smaller variable costs. On-demand pricing is designed for proof of concepts or low commitment utilization. This gives developers the option to shut down projects instantly or as needed.

3 Year Reserved Instance Pricing

Region: US East (N. Virginia) ▾		
3-Year Reserved Instances		
DW Node Class (Reserved)	Upfront	Hourly
dw1.xlarge	\$3,000	\$0.114 per Hour
dw1.8xlarge	\$24,000	\$0.912 per Hour
dw2.large	\$1,325	\$0.050 per Hour
dw2.8xlarge	\$21,200	\$0.800 per Hour

source: “Amazon Redshift Pricing,” 2015

Reserved Instances offers a 75% discount in pricing compared to on-demand.

Additionally, it asks for a low, one-time payment to reserve each instance and in turn receive a significant discount on the hourly charge for that instance. There are three Reserved Instance types (Light, Medium, and Heavy Utilization Reserved Instances) that enable you to balance the amount you pay upfront with your effective hourly price.

When comparing on-demand vs reserved instances by the TB, the difference between the two are substantial. For example, the Oracle 30 TB database would compare as follows:

Estimated Price for 30 TB per Year			
	On-Demand	1yr RI	3yr RI
dw1.xlarge (2 TB HDD)	\$111,690.00	\$65,760.00	\$29,970.00
dw1.8xlarge (16TB HDD)	\$111,690.00	\$65,760.00	\$29,970.00
dw2.large (0.16TB SSD)	\$410,640.00	\$263,820.00	\$164,940.00
dw2.8xlarge (2.56 TB SSD)	\$492,750.00	\$330,540.00	\$164,940.00

7/20/2015

INL-880 - Capstone Proposal: McGinley & Etter -Final Draft

source: appendix, INL 880: Capstone Product Worksheet

These costs are factored based off of three tiers; compute node hours, backup storage, and data transfer.

Compute node hours are the total hours that are run against all of the compute nodes per billing period (which is typically monthly). Compute nodes are billed 1 unit per node per hour. For example, let's assume running a persistent run for a single (read: one) node would be approximately 720 hours. The instance hours billed would be 720. Additionally, Amazon will not charge for the leader nodes that are automatically created. So if you have two nodes (with one or more leader nodes) running persistently, you will be billed for 1,440 instance hours (read: 2 nodes running for 720 hours).

Backup storage is the additional manual snapshot of the data warehouse that is desired. To note, Amazon will not charge for storage up to 100% of the provisioned storage of an active warehouse cluster. For example, it is estimated that if two active nodes are provisioned to equal approximately 30TB of storage, Amazon will provide 30TB of backup storage for no additional cost.

The actual combined annual cost (on-demand vs reserved instance) using Amazon's calculator is as follows:

Actual Calculation (all at 100% utilization) 30TB per year (or as close as possible) *Amazon does not include support costs in initial estimation as it is listed below				
Node Type	Nodes Required	On-Demand w/ Support	1yr RI w/ Support	3yr RI w/ Support
dw1.xlarge (2 TB HDD)	15	\$123,195.60	\$71,515.42	\$64,935.61
dw1.8xlarge (16TB HDD)	2	\$131,408.64	\$76,263.14	\$69,244.35

dw2.large (0.16TB SSD)	188	\$442,939.92	\$285,567.86	\$353,817.91
dw2.8xlarge (2.56 TB SSD)	12	\$544,975.56	\$364,730.16	\$361,199.14

source: appendix, INL 880: Capstone Product Worksheet

One factor that was not accounted for until the actual calculations were in place was the additional support charge from Amazon:

	Basic	Developer	Business	Enterprise
Pricing	Included	\$49/month	Greater of \$100 - or - 10% of monthly AWS usage for the first \$0–\$10K 7% of monthly AWS usage from \$10K–\$80K 5% of monthly AWS usage from \$80K–\$250K 3% of monthly AWS usage over \$250K	Greater of \$15,000 - or - 10% of monthly AWS usage for the first \$0–\$150K 7% of monthly AWS usage from \$150K–\$500K 5% of monthly AWS usage from \$500K–\$1M 3% of monthly AWS usage over \$1M

source: “AWS Support Pricing” 2015

It is estimated that only three months of application development are required from three developers to configure the system. Additionally, once provisioned, one full time DBA would be sufficient for ongoing administration and configuration. This and all following resource assumptions do not include any additional education or training costs.

	Redshift			
Upfront Resource Work	Assume baseline of 3 man-months of application development of application development (Developer salary of \$120,000.00 per year) 3 Total Developers	\$90,000.00		
			\$90,000.00	<-- Year 1
Ongoing Resource Work				
	Assumes 1 full time DBA's with salaries of \$120,000 per year	\$120,000.00	\$120,000.00	<-- Year 2

source: appendix, INL 880: Capstone Product Worksheet

Cost Structure and Resources Needed

Redshift	On-Demand Instance 100% utilization. 30TB per year (or as close as possible). Business Support. dw.xlarge 2 TB HDD. 15 nodes.	1yr Reserved Instance 100% utilization. 30TB per year (or as close as possible). Business Support. dw.xlarge 2 TB HDD. 15 nodes.	3yr Reserved Instance 100% utilization. 30TB per year (or as close as possible). Business Support. dw.xlarge 2 TB HDD. 15 nodes.
Total Upfront Resource Costs	\$90,000.00	\$90,000.00	\$90,000.00
Total Ongoing Resource Costs	\$120,000.00	\$120,000.00	\$120,000.00
Total Upfront Infrastructure Cost	N/A	\$40,354.18	\$48,412.45
Total Ongoing Infrastructure Cost	\$123,195.60	\$31,161.24	\$16,523.16
Total Year 1	\$213,195.60	\$161,515.42	\$154,935.61
Total Year 2	\$243,195.60	\$151,161.24	\$136,523.16
Total Year 3	\$243,195.60	\$151,161.24	\$136,523.16
3 Year TCO	\$699,586.80	\$463,837.90	\$427,981.93

source: appendix, INL 880: Capstone Product Worksheet

Google BigQuery

BigQuery is Google Cloud Engine's offering for an online analytical processing (OLAP) database. Specifically BigQuery is Google opening its infrastructure, specifically the Dremel Analytics Processing Architecture, for commodity usage. Dremel was designed by Google for ad hoc analytics queries. Google's description is "Dremel is a query service that allows you to run SQL-like queries against very, very large data sets and get accurate results in mere seconds. You just need a basic knowledge of SQL to query extremely large datasets in an ad hoc manner.

At Google, engineers and non-engineers alike, including analysts, tech support staff and technical account managers, use this technology many times a day." (Sato, 2). The speed of BigQuery and Dremel's queries are extremely fast, often faster than even Amazon's Redshift offering and even easier to maintain. The BigQuery offering from Google offers Zero operational maintenance which, unlike Redshift, does not require users to monitor and

understand the compute and data needs. BigQuery automatically resizes clusters as needed

without any human interaction. Google accomplishes this by leveraging all of their existing infrastructure built for their products such as AdWords, Search, Books and many others, which was specifically designed for elasticity and scalability. As data grows, BigQuery will easily be able to ingest it and make it consumable for analytics usage. It should be made very clear that BigQuery is designed specifically for analytics and not transactional usage, because of this BigQuery does not allow for future updating of records. Once data enters BigQuery it is immutable data, and can really only be updated by completely removing the records and entering new ones. As Google explains for BigQuery is not for all use cases but for the cases it is “ By using Dremel instead of MapReduce on about two-thirds of all my analytic tasks, I was able to finish the job by lunch time. And if you’ve ever eaten lunch at Google, you know that’s a big deal.” (Sato, 7)

SQL Compliance

BigQuery supports many of the ANSI SQL standards. Specifically, according to Google’s reference manual it supports “SELECT, WITHIN, FROM, FLATTEN, JOIN, WHERE, GROUP BY, HAVING, ORDER BY, and LIMIT” (BigQuery: Query Reference, 2015). This is enough support of ANSI SQL to cover all of the standard E-Commerce analytics queries and to also provide us with the ability to do many ad-hoc queries. Notice that the INSERT and UPDATE functions are not supported by BigQuery. Unlike Redshift, BigQuery does not allow any individual insert of data. Updating data is also explicitly not allowed and must be accomplished through dropping data and then reinserting it. Insertion of data is done entirely programmatically and the SQL specification that BigQuery uses supports this. The one interesting part of ANSI SQL that BigQuery does not support is the wild card (*) syntax for

choosing columns when running a SELECT query. Instead of using the wildcard in a query the inspector of data must specify all the columns they are interested in explicitly. Much like Redshift new tables in BigQuery must be described with a schema, before data can be inserted. This schema helps BigQuery optimize the storage of data across the cluster.

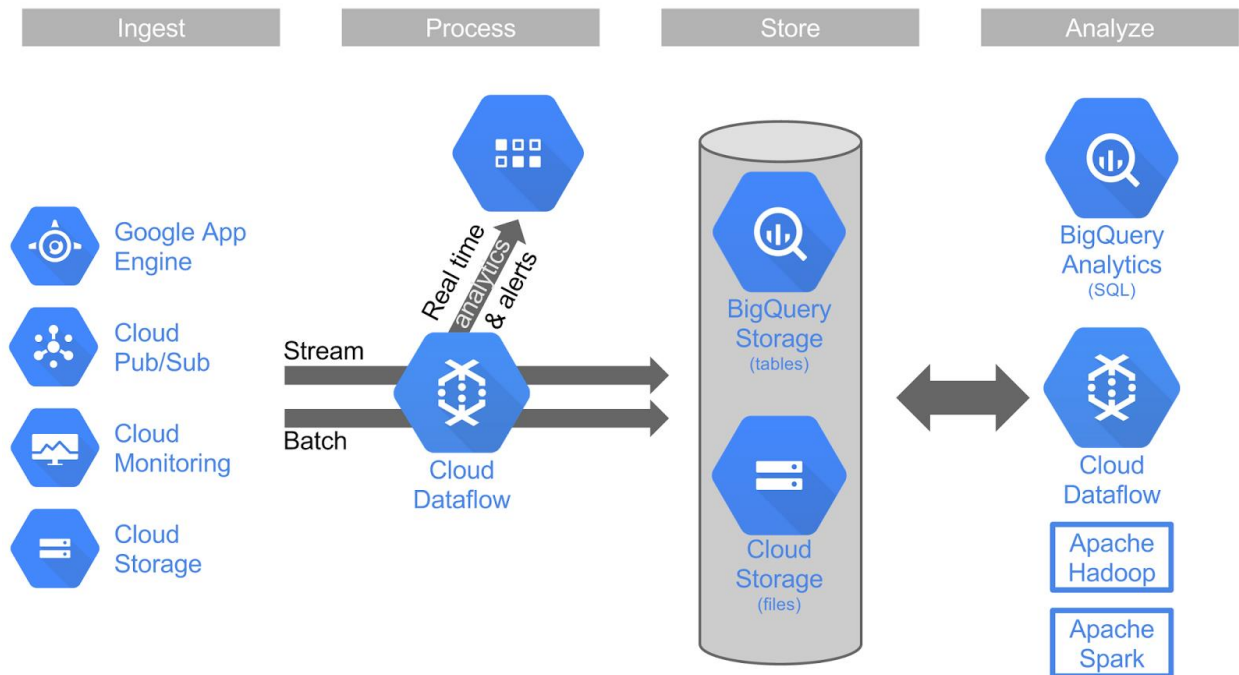
Performance and Scalability

Google, unlike many other cloud providers, runs on it's own own global fiber network. When every millisecond of latency counts, Google ensures that content is delivered quickly. Google states that streaming data through can query 100,000 rows per second to enable real-time analysis of data. ("Why Google Cloud Platform," 2015)

BigQuery was designed to scale to the exabyte size of data. This is a problem that not many companies in the world will face. Google however does face this problem with it's search engine, and BigQuery, known internally as Dremel, is the answer. Dremel has allowed Google to analyze petabytes of data for trends or answers in seconds. Dremel, and therefor BigQuery are capable of performance, much better than MapReduce(MR), "MR gains an order of magnitude in efficiency by switching from record-oriented to columnar storage (from hours to minutes). Another order of magnitude is achieved by using Dremel (going from minutes to seconds)." (Melnik. 2010). Dremel is the evolution of all the work Google did to create MapReduce and allows large datasets to be quickly explored for new trends, BigQuery promises to give the exact same capabilities to any consumers willing to utilize the service.

Integrations

Google offers a wide array of tools that integrate with BigQuery in order to streamline the experience:

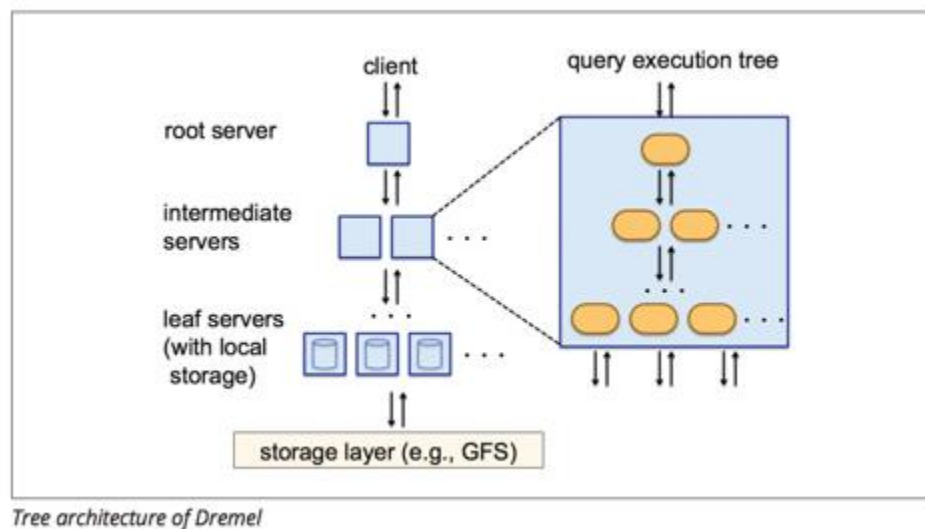


source: Vambenepe. April 16, 2015

Google provides multiple out-of-the-box solutions and continues to develop new technologies to improve the user experience. Additionally, BigQuery integrates with the most commonly used visualization, business intelligence, and ETL tools. This includes Tableau and snapLogic which utilize the “open APIs provided by Google Cloud Storage and BigQuery.” (BigQuery. “Third-party Tools and Services”)

Architecture

BigQuery utilizes the advent of the commodity server to maximize its performance and storage. Instead of storing as much data as possible on one server BigQuery distributes it across many different servers, called leafs. BigQuery utilizes columnar data storage and a tree structure to optimize query time.



source: (Sato. 2012)

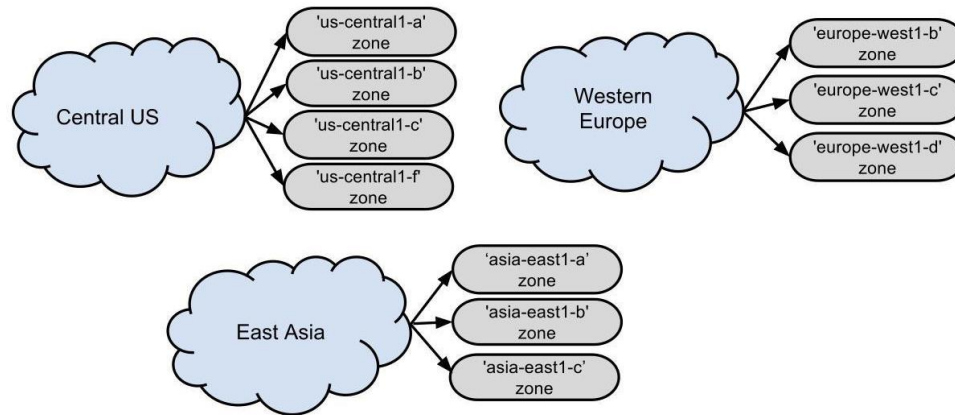
This tree structure is exactly like the tree structures studied in computer science. There is a root server that analyzes the initial query and creates an optimized c++ program to be run across all of the leaf nodes on the tree. The leaf nodes run the query and yield their results to the nodes above them who join the various answers together and eventually reaching the root node which performs the final join of data and returns the result to the user. Google has also designed BigQuery so that the more data you add, nodes are added to the tree seamlessly and without the need for interaction. Redshift requires manual intervention to resize a cluster. This is

a huge win for the end user as it reduces the cost of operations and knowledge or talent required to maintain the data store. In the end BigQuery is a data store that could even be used and maintained by a relatively smart analytics user, without any real technical knowledge about server management or lots of programming experience. (Stato. 2012)

Similar to Redshift all of BigQuery's backup and disaster recovery is maintained as part of Google's infrastructure. Google does a significant amount of replication of data, into slaves and cold storage. However Google's infrastructure takes this one step further and seamlessly and automatically replaces failing hardware in the backend. This removes backups and the replacement of a failing virtual server not a factor for consideration. In short, for BigQuery to have a catastrophic failure would mean that Google would also have to have a catastrophic failure.

Security

As of April 2015, BigQuery has expanded its capabilities into European zones. This contributes significantly to the scalability and redundancy options in term of flexibility. Users are now able to distribute resources across multiple zones, or isolated location within a region, in multiple regions, or collections of zones with high-bandwidth and low-latency connected to each other.



source: “Google Cloud Platform Regions & Zones,” 2015

Taking advantage of the zones available will protect users from unplanned downtime or failure.

Google has opened its doors to an entire ecosystem of enterprise applications for BigQuery by adding data expiration controls along with row-level permissions. “Row-level permissions eliminate the need to create different views for different users, allowing secure shared access to systems such as finance or HR. This ensures that you get the information that’s relevant to you. In addition, data in BigQuery will be encrypted at rest.” (“Google Cloud Platform Regions & Zones,” 2015)

BigQuery uses Access Control Lists (ACLs) to manage permissions on projects and datasets. Further, ACLs are not directly supported on table as a table will inherit its ACL from the dataset that contains it. (“BigQuery Access Control,” 2015) Google’s Cloud Platform shares the same infrastructure with Google Apps. The security, compliance, rigorous audit trail and certification efforts are substantial. Google is constantly undertaking the required tasks to be approved and accredited by the most popular third party audits for data safety, privacy, and security:

Audit/ Program	Explanation
ISO 27001	One of the most widely recognized, internationally accepted independent security standards, and Google received the certification for Google Cloud Platform.
SOC2,, SOC3 public audit report, and ISAE 3402	Google has successfully completed the SOC2, SSAE 16 Type II audit, and its international counterpart ISAE 3402 Type II audit, to document and verify the data protections in place for their services.
HIPAA	In 2014, Google started entering into Business Associate Agreements (BAAs) to allow Google Apps customers to support HIPAA regulated data.
FISMA Moderate accreditation for Google App Engine	
Payment Card Industry data (PCI DSS v3.0)	Google Cloud Platform has been validated for compliance with the Payment Card Industry (PCI) Data Security Standards.
US Department of Commerce Safe Harbor Program	Google will remain enrolled in this program or another replacement program (or will adopt a compliance solution which achieves compliance with the terms of Article 25 of Directive 95/46/EC)
SAS70 and SSAE16	Google is certified for SAS70 and SSAE16 which makes it simpler for organizations to go through certification. Companies must only certify from the path from source code to the App Engine platform.

source: “Google Security Whitepaper,” 2015 & “Total Economic Impact of Google Cloud Platform,” 2014

The expected work to design and implement the PEGGY system through the use of BigQuery is as follows:

	Big Query			
Upfront Resource Work	Assume baseline of 3 man-months of application development of application development (Developer salary of \$120,000.00 per year)	\$30,000.00	\$30,000.00	<-- Year 1
	Assumed no DBA required			
Ongoing Resource Work	1 Developer for ongoing support at \$135,000 per year	\$135,000.00	\$135,000.00	<-- Year 2

source: appendix, INL 880: Capstone Product Worksheet

The expected requirements to design and implement the PEGGY system through the use of BigQuery is as follows:

BigQuery
Data Warehouse
Storage 3,000 GB (30TB)
Streaming Inserts 0 rows

7/20/2015

INL-880 - Capstone Proposal: McGinley & Etter -Final Draft

Interactive Queries 30 TB
Batch Queries 0 TB
\$205.00
Monthly total: \$205.00

source: appendix, INL 880: Capstone Product Worksheet

The expected combined total costs for the work involved to design and implement the PEGGY system through the use of BigQuery is as follows:

BigQuery	Storage 3,000 GB (30TB) Streaming Inserts 0 rows Interactive Queries 30 TB Batch Queries 0 TB
Total Upfront Resource Costs	\$30,000.00
Total Ongoing Resource Costs	\$135,000.00
Total Upfront Infrastructure Cost	N/A
Total Ongoing Infrastructure Cost	\$2,460.00
Total Year 1	\$32,460.00
Total Year 2	\$137,460.00
Total Year 3	\$137,460.00
3 Year TCO	\$307,380.00

source: appendix, INL 880: Capstone Product Worksheet

Hadoop

Hadoop is a programming framework designed for the storage and analysis of large volumes of data. Specifically The Apache Foundation describes Hadoop as:

“a framework that allows for the distributed processing of large data sets across clusters of computers using simple programming models. It is designed to scale up from single servers to thousands of machines, each offering local computation and storage. Rather than rely on hardware to deliver high-availability, the library itself is designed to detect and handle failures at the application layer, so delivering a highly-available service on top of a cluster of computers, each of which may be prone to failures.” (“Welcome to Apache,” 2014.)

The Hadoop framework's primary algorithm of Map/Reduce was inspired by Google who used the algorithm for years for calculating the value of individual website pages in their Pagerank Algorithm (Dean. 2004) Google has long since abandoned the Map/Reduce implementation of Pagerank in favor newer more efficient algorithms. Despite, this change by Google many companies still find extensive use from the Hadoop implementation of Map/Reduce which also provides tools for organizing groups of servers, and storage of data in a manner that can be most effectively utilized by the algorithm. The data warehouse part of the Hadoop is known as the Hadoop Distributed File System (HDFS). HDFS provides redundancy, and other safeguards to ensure that data is not easily lost, and also provides the mechanism for moving data so that it has a locality related to the machines processing it. HDFS' primary feature is that it stores unstructured data. As long as what is being stored in HDFS is a file it can be stored for later analysis. Most often you will see data in the form of Javascript Object Notation, Column Separated Files, or Extensible Markup Language. ("JaqlOverview," 2014)

SQL Compliance

Hadoop and the Map/Reduce algorithm provide no SQL compliance by default. However the Apache Foundation, took over a project from Facebook, who preferred to report on their data using a SQL like interface, called Hive. Hive provides a SQL like interface on top of Map/Reduce and semistructured data. By utilizing a Serializer/Deserializer per file type Hive is able to translate SQL queries into Map/Reduce jobs that act on files. Specifically the Apache foundation states

“The Apache Hive™ data warehouse software facilitates querying and managing large datasets residing in distributed storage. Hive provides a mechanism to

project structure onto this data and query the data using a SQL-like language called HiveQL. At the same time this language also allows traditional map/reduce programmers to plug in their custom mappers and reducers when it is inconvenient or inefficient to express this logic in HiveQL.” (Apache Hive™, 2014)

The Hive tool and HiveQL allow traditional database analysts to use a tool that is extremely familiar to them while analyzing data at scale. Despite all of the group's attempts to make HiveQL as much like SQL as possible there are some noticeable differences with ANSI SQL. Specific deviations from ANSI SQL focus around the operations of JOIN, GROUP BY, and COLLECT SET. In HiveQL JOIN operations “Only equality joins, outer joins, and left semi joins” and no other types are supported “ as it is very difficult to express such conditions as a map/reduce job.” (“LanguageManualJoins,” 2014) In HQL all columns in a select must be present in the GROUP BY or accessed only through an aggregation function, so that no guessing work is required to determine what else should be shown in the results. HQL enforces this by making access of columns outside of a GROUP BY or aggregation function invalid SQL. Collect set is a special function that “allows a column not used in the group by to be aggregated into a set. The values in the set are accessible using normal array-like syntax and can be used the same way as any column in the original table.”(Spry, 2013)

Performance and Scalability

Hadoop is used as the storage and analysis backend for companies that store large amounts of data. Facebook, a well known advocate of Hadoop, released the first open source version of Hadoop. (Borthakur . 2010) Facebook has used Hadoop into the high Petabyte size of

data, across their social graph data storage. The tools created for Hadoop, such as Parquet, which provides Columnar data storage for Hadoop, are also quickly increasing performance capabilities. Compared to Amazon Redshift, Hadoop can scale and perform just as well however it requires much more engineering knowledge to do so. Hadoop clusters do not automatically scale and new technologies are not automatically integrated in a maintenance cycle. The development team must read and understand how all new tools for Hadoop fit in the infrastructure, whether they are useful for the team and then come up with a plan to easily deploy the additions. For simple MapReduce jobs on data you have stored in Amazon's Simple Storage Service you can utilize Elastic MapReduce (EMR). This does not maintain an in memory HDFS cluster, instead it starts up a cluster of computers and loads data from S3 and immediately performs the analysis on the data, returns the results and destroys the cluster. EMR is also completely compatible with Hive. This allows companies to cheaply store all of their data in S3, and without having the need for experts to maintain a Hadoop cluster. The trade off for this ease of use of MapReduce is that there is very bad data locality and you pay on cluster startup from loading all data across Amazon's network. This slows down analysis and will not be as fast as maintaining a single cluster. So Amazon's Hadoop provides many tradeoffs for users to decide what is best for them.

Integrations

Hadoop cannot be easily integrated into all existing Business Intelligence tools because of its lack of full ANSI SQL compliance. Because of this, it is more likely that resources will be required to invest in new Business Intelligence tools for analysis of the data. Some of the tools that can work with Hadoop include Tableau and Terradata, both of which also try to bridge the

gaps that exist between traditional Business Intelligence tools and the programming paradigms associated with the Hadoop ecosystem. There are also several companies providing Enterprise Grade Hadoop Platforms on Amazon. These companies include Hortonworks, ParAccel and Cloudera. *Integrating Hadoop Into Business Intelligence and Data Warehousing* explains particularly how Hortonworks is working hard to bridge the gap associated with Hadoop

Hortonworks focuses on innovating the core of open source Apache Hadoop in ways that make Hadoop enterprise grade and therefore more applicable to more user organizations. Hortonworks' strategy is to distribute 100% open source Apache Hadoop, with additional operational, data, and platform services from the open source community, all packaged as the Hortonworks Data Platform (HDP). Multi-tenancy is built into HDP, so it can be a shared enterprise infrastructure instead of a silo, and HDP 1.2 beefs up security, which is the leading concern of Hadoop users. Hortonworks is a major contributor to open source Hadoop technologies, and it has recently shown leadership in the design of Apache HCatalog (metadata services for the Hadoop ecosystem), Apache Ambari (cluster management and monitoring for HDFS), and high availability for NameNode in Hadoop 2.0. (Russom, p. 30)

Over time these tools will continue to improve as a large number of companies, including IBM are working hard to improve them, and contributing back to the open source project. This will give us continual improvements, faster than what one organization can provide.

Architecture

Hadoop has a cluster based architecture composed of clusters of commodity machines. These servers can also be virtual machines provided by Amazon Web Services, Google Compute Engine, Digital Ocean or others. The primary data storage mechanism of Hadoop is HDFS, which stores data across many nodes and also stores the data much like a Redundant Array of Inexpensive Disks (RAID). The workers' nodes utilize technologies such as MapReduce, Hadoop, Spark and other similar tools. For managing the actual MapReduce jobs, Hadoop uses Yet Another Resource Negotiator (YARN).

YARN remedies the scalability shortcomings of “classic” MapReduce by splitting the responsibilities of the jobtracker into separate entities. The jobtracker takes care of both job scheduling (matching tasks with tasktrackers) and task progress monitoring (keeping track of tasks, restarting failed or slow tasks, and doing task bookkeeping, such as maintaining counter totals). (White, 170)

The layers of clusters provide operational redundancy to Hadoop and are part of the disaster recovery process. Node replacement is still a manual task that requires knowledge about what to do, however this can also eventually be scripted. The final piece of the disaster recovery plans is to backup data to a cloud service such as Amazon's S3 or Google's Cloud Storage. Backing up to these services will allow you to rebuild a cluster from scratch should it fail, catastrophically.

Security

Hadoop nodes have the standard security provided by all of the Cloud Services. The cold storage has it's own permissions structures, and servers have access controls provisioned to them

by operations. After the basic security policies are set, an analysis of Hadoop's built in security is required. Hadoop uses Kerberos for security and only has a few authorization settings.

Although Hadoop can be configured to perform authorization based on user and group permissions and Access Control Lists (ACLs), this may not be enough for every organization. Many organizations use flexible and dynamic access control policies based on XACML and Attribute-Based Access Control. Although it is certainly possible to perform these level of authorization filters using Accumulo, Hadoop's authorization credentials are limited. (Smith, 2013)

The other issue with Kerberos is that it's difficult to setup with extremely fine grained permissions that need to be set right or they will conflict and deny permission. The other issue with Hadoop, lies in HDFS. Data written to disk is not encrypted at all, this means whoever gains access to the servers with data on them can read the data without requiring further authorization.

Currently, data is not encrypted at rest on HDFS. For organizations with strict security requirements related to the encryption of their data in Hadoop clusters, they are forced to use third-party tools for implementing HDFS disk-level encryption, or security-enhanced Hadoop distributions. (Smith, 2013)

This means manual development resources will be required to enhance HDFS's security, or engage with a vendor who has a security focused version of Hadoop to sell, completely changing the cost structure. The final option to avoid the data at rest security issue is to not store anything of sensitivity on the Hadoop cluster. This option will affect what can be done with the data so it is expected that in order to move forward, the server level security must be acceptable enough to continue.

Cost Structure

The capability to run Hadoop locally such as an on-premise server farm or on any of the cloud platforms from Google Compute Engine to Amazon Web Services will impact the cost structure. Determining which commodity servers to use would also be based on the talent available to build the systems and maintain them. This shifts the majority of the costs to the talent needed to hire or develop so that the Hadoop cluster system can be properly maintained. Initially there will be some issues with a skills shortage much like Forrester suggests in this quote from Progressive Digital Media Technology News. “The shortage of Hadoop skills will quickly disappear as enterprises turn to their existing application development teams to implement projects such as managing data lakes and developing MapReduce jobs using Java, according to” (10 hadoop predictions for 2015. 2014). To this extent, an additional invest in training the internal development team through e learning platforms such as Cloudera or Hortonworks to further develop the talent to maintain Hadoop inhouse. Specifically the changes will greatly affect the current Database Administrators as they will need to be cross trained. Further training costs would expand to all of the development staff on at least the basics. Based on the public calculators available, the following estimated costs have been determined.

	Hadoop Resources			
Upfront Resource Work	Assume baseline of 3 man-months of application development of application development (Developer salary of \$120,000.00 per year)	\$30,000.00	\$30,000.00	Year 1 Cost
Ongoing Resource Work	1 Hadoop expert for ongoing support at \$145,000 per year	\$145,000.00	\$145,000.00	Year 2 Cost

source: Appendix, INL 880: Capstone Product Worksheet

As mentioned earlier during the cost analysis review of Amazon Redshift, the resource costs for training and education have been omitted from this study. However, it should be noted

that Hadoop, more so than Redshift or BigQuery, has been identified to have the largest learning requirement.

Compute Engine
21,900 total hours per month
Instance type: n1-standard-8-preemptible
Region: United States
Total Estimated Cost: \$1,839.60
(Sustained Use Discount: 30%)
Sustained Usage Discount Monthly Breakdown:
1st ¼ - 5,475.0 hrs @ 0.0% off: \$657.00
2nd ¼ - 5,475.0 hrs @ 20.0% off: \$525.60 (\$131.40 saved)
3rd ¼ - 5,475.0 hrs @ 40.0% off: \$394.20 (\$262.80 saved)
4th ¼ - 5,475.0 hrs @ 60.0% off: \$262.80 (\$394.20 saved)
(Effective Hourly Rate: \$0.084)
Persistent Disk
SSD storage: 0 GB
Storage: 3,000 GB
Snapshot storage: 0 GB
\$120.00
Monthly total: \$1,959.60

source: Appendix, INL 880: Capstone Product Worksheet

	Compute Engine 21,900 total hours per month Instance type: n1-standard-8-preemptible Region: United States Persistent Disk Storage: 3,000 GB (ie 30TB)
Hadoop	
Total Upfront Resource Costs	\$30,000.00
Total Ongoing Resource Costs	\$145,000.00
Total Upfront Infrastructure Cost	N/A
Total Ongoing Infrastructure Cost	\$23,515.20
Total Year 1	\$53,515.20
Total Year 2	\$168,515.20
Total Year 3	\$168,515.20
3 Year TCO	\$390,545.60

source: Appendix, INL 880: Capstone Product Worksheet

Impact on resources

As mentioned in the cost analysis, the Hadoop infrastructure provides all new ways to think about Data Analysis. There is a plethora of tools being built to run on top of Hadoop for data analysis. Training will help get engineering and data analysis teams up to speed on the basics of the Hadoop Framework, MapReduce and Hive. The new tools being developed for Hadoop will require continuous evaluation to ensure they fit with the types of analysis and budget desired. This will require the developers to do a significantly larger amount of research before integrating the new technology. This will require a familiarity with engaging Open Source communities in order to obtain answers to questions. This is much different than engaging a support system backed by the service level agreements offered by Oracle, Google, and Amazon. In the end there is a lot of cost in the situation and new talent would be needed to successfully deploy Hadoop based system. There would also have significantly more overhead from an operations perspective compared to Redshift or Hadoop.

Comparing the Cloud Tools

Miles Ward, former Senior Manager of Solutions Architecture at Amazon Web Services and now Global Head of Solutions at Google, has been running a series of blog posts entitled, Understanding Cloud Pricing. In June of 2015, he put together an example cluster comparing BigQuery, Amazon Redshift, and Apache Hadoop. He insisted that although Hadoop is an open source query engine, the combination of the storage capabilities and proprietary solutions available for analysis are equivalent to the cloud-based data warehousing solutions BigQuery and Redshift. His goal throughout his study was to compare and analyze three systems (Hadoop,

Redshift and BigQuery) capabilities surrounding the storage of massive amounts of data and analytical reporting running exclusively on the public cloud. (Ward. June 19, 2015)

Ward's study uses the following parameters for the tools assessed.

BigQuery	Hadoop	Redshift On-Demand	Redshift 1 yr Reserved	Redshift 3yr Reserved
100 users, 40 queries each per day, with 100 GB average query size. (ie 4000 queries per day w/ 12,000 TB data analyzed per month (simplifying to 30 day month).	Compute Engine 73,000 total hours per month Instance type: n1-highmem-16 Region: United States Persistent Disk Storage: 1,000,000 GB	On-Demand 63 nodes. dw1.8xlarge - 16TB HDD (1008TB) 100% utilized. Free Support	1yr Reserved. 63 nodes. dw1.8xlarge - 16TB HDD (1008TB) 100% utilized. Free Support	3yr reserved 63 nodes. dw1.8xlarge - 16TB HDD (1008TB) 100% utilized. Free Support

source: appendix, INL 880: Capstone Produc Worksheet

The calculations uncovered that the on-demand pricing for BigQuery and Hadoop outmatch the cost of Redshift by approximately 78%. To note, this comparison, although produced by a Google employee, utilized the publicly available calculators to factor costs.

	BigQuery	Hadoop	Redshift On-Demand
Total Upfront Infrastructure Cost	N/A	N/A	N/A
Total Ongoing Infrastructure Cost	\$959,994.00	\$1,206,028.80	\$3,787,065.60
Total Year 1	\$959,994.00	\$1,206,028.80	\$3,787,065.60
Total Year 2	\$959,994.00	\$1,206,028.80	\$3,787,065.60
Total Year 3	\$959,994.00	\$1,206,028.80	\$3,787,065.60
3 Year TCO	\$2,879,982.00	\$3,618,086.40	\$11,361,196.80

source: appendix, INL 880: Capstone Produc Worksheet

Additionally, if the organization was to commit to a three year reserved instance of Amazon Redshift, there would need to be an upfront payout of approximately \$1.5 million in addition to the monthly ongoing cost of \$42,058.17. Google offers a sustained usage discount for running on Google's Compute Engine. This is an automatic discount that Google offers for virtual machines which further reduces the overall cost and eliminates the need for additional license fees. (Ward. June 19, 2015)

	BigQuery	Hadoop	Redshift On-Demand	Redshift 1 yr Reserved	Redshift 3yr Reserved
Total Upfront	N/A	N/A	N/A	\$1,260,000.00	\$1,512,000.00

Infrastructure Cost					
Total Ongoing Infrastructure Cost	\$959,994.00	\$1,206,028.80	\$3,787,065.60	\$79,319.52	\$504,698.04
Total Year 1	\$959,994.00	\$1,206,028.80	\$3,787,065.60	\$1,339,319.52	\$2,016,698.04

source: appendix, INL 880: Capstone Produc Worksheet

However, considering the data warehouse has a shelf life of three years, the total cost of ownership combined with Amazon's three year reserved instance option results in a slightly more level playing field when comparing Ward's findings. In fact, over three years, Redshift's total cost of ownership is lower than the cost of Hadoop making Redshift and BigQuery the two least expensive options.

	BigQuery	Hadoop	Redshift On-Demand	Redshift 1 yr Reserved	Redshift 3yr Reserved
Total Upfront Infrastructure Cost	N/A	N/A	N/A	\$1,260,000.00	\$1,512,000.00
Total Ongoing Infrastructure Cost	\$959,994.00	\$1,206,028.80	\$3,787,065.60	\$79,319.52	\$504,698.04
Total Year 1	\$959,994.00	\$1,206,028.80	\$3,787,065.60	\$1,339,319.52	\$2,016,698.04
Total Year 2	\$959,994.00	\$1,206,028.80	\$3,787,065.60	\$79,319.52	\$504,698.04
Total Year 3	\$959,994.00	\$1,206,028.80	\$3,787,065.60	\$79,319.52	\$504,698.04
3 Year TCO	\$2,879,982.00	\$3,618,086.40	\$11,361,196.80	\$1,497,958.56	\$3,026,094.12

source: appendix, INL 880: Capstone Produc Worksheet

Amazon offers multiple cost and processing structures for producing Redshift clusters. For the purpose of the project requirements, even the best pricing available positioned Redshift as the most expensive cloud solution when compared to BigQuery and Hadoop.

Oracle	BigQuery	Hadoop	Redshift On-Demand	Redshift 1yr Reserved	Redshift 3yr Reserved
Oracle Database Subscriber Edition & Oracle Real Application Cluster (RAC) Server Hardware: 30 Servers (8 cores/server) w/ 32 GB RAM Storage Hardware: 30 TB SAN (usable)	Storage 3,000 GB (30TB) Streaming Inserts 0 rows Interactive Queries 30 TB Batch Queries 0 TB	Compute Engine 21,900 total hours per month Instance type: n1- standard-8- preemptible Region: United States Persistent Disk Storage: 3,000 GB (ie 30TB)	On-Demand Instance 100% utilization. 30TB per year (or as close as possible). Business Support. dw.xlarge 2 TB HDD. 15 nodes.	1yr Reserved Instance 100% utilization. 30TB per year (or as close as possible). Business Support. dw.xlarge 2 TB HDD. 15 nodes.	3yr Reserved Instance 100% utilization. 30TB per year (or as close as possible). Business Support. dw.xlarge 2 TB HDD. 15 nodes.
\$4,850,000.00	N/A	N/A	N/A	\$40,354.18	\$48,412.45
\$992,600.00	\$2,460.00	\$23,515.20	\$123,195.60	\$31,161.24	\$16,523.16
\$4,850,000.00	\$2,460.00	\$23,515.20	\$123,195.60	\$71,515.42	\$64,935.61
\$992,600.00	\$2,460.00	\$23,515.20	\$123,195.60	\$31,161.24	\$16,523.16
\$992,600.00	\$2,460.00	\$23,515.20	\$123,195.60	\$31,161.24	\$16,523.16
\$6,835,200.00	\$7,380.00	\$70,545.60	\$369,586.80	\$133,837.90	\$97,981.93

source: appendix, INL 880: Capstone Produc Worksheet

Further, once the assumed resource calculations were included, the total cost of ownership remained in favor of BigQuery as it produces \$83,165.60 in savings over three years when compared to Hadoop. The savings are \$120,601.93 when compared to Redshift's three year reserved instance. Most striking was the \$8,387,820.00 in savings when comparing to the on-premise Oracle upgrade.

	Oracle	BigQuery	Hadoop	Redshift On-Demand	Redshift 1yr Reserved	Redshift 3yr Reserved
Total Upfront Resource Costs	\$780,000.00	\$30,000.00	\$30,000.00	\$90,000.00	\$90,000.00	\$90,000.00
Total Ongoing Resource Costs	\$540,000.00	\$135,000.00	\$145,000.00	\$120,000.00	\$120,000.00	\$120,000.00
Total Upfront Infrastructure Cost	\$4,850,000.00	N/A	N/A	N/A	\$40,354.18	\$48,412.45
Total Ongoing Infrastructure Cost	\$992,600.00	\$2,460.00	\$23,515.20	\$123,195.60	\$31,161.24	\$16,523.16
Total Year 1	\$5,630,000.00	\$32,460.00	\$53,515.20	\$213,195.60	\$161,515.42	\$154,935.61
Total Year 2	\$1,532,600.00	\$137,460.00	\$168,515.20	\$243,195.60	\$151,161.24	\$136,523.16
Total Year 3	\$1,532,600.00	\$137,460.00	\$168,515.20	\$243,195.60	\$151,161.24	\$136,523.16
3 Year TCO	\$8,695,200.00	\$307,380.00	\$390,545.60	\$699,586.80	\$463,837.90	\$427,981.93

source: appendix, INL 880: Capstone Produc Worksheet.

Conclusion and Recommendation

A high level overview has been performed to determine the characteristics, intricacies, and overall costs of the various platforms. The project focus has mostly been on the cost of goods and also the cost in skill acquirement. If there was an opportunity to continue the research, the next step would be to take an example dataset and analyze performance variations between Redshift, BigQuery, Hadoop, and Oracle. Secondly, it would also benefit the findings by examining how table structure may affect the performance of these tools. For example, how would performance and cost be predicted when comparing BigQuery's automatic scaling versus Redshift's manual administration and configuration. Would faster processing power outweigh a predictable cost structure -would it matter? All of this additional analysis would allow us to predict the required skillset of hired talent, workload to maintain a performant datastore, and the true cost of operating a cloud based data warehousing solution.

The initial assumption was that Amazon Redshift, a cost completely customizable, secure, competitive, flexible, SQL compliant, cloud-based data warehouse as a service would be the leader. However, based on the research, the significant cost savings, predicted performance and scalability speeds, security, minimal skillset, and tolerable compliance to SQL casts Google's BigQuery as being the tool of choice. This was especially unexpected as BigQuery offers a limited SQL experience. Yet, this customized language is believed to require minimal learning to effectively utilize the tool. This is a significant change from the original project scope.

During the early stages of research, it was discovered that Amazon's RDS service, apart from being offered as an on-demand or reserved instance service managed by Amazon, is too closely matched to the legacy on-premise Oracle system. The similarities were too parallel to the

legacy system and, apart from the concern of a possible new schema, it was decided to remove RDS and replace it with Google's BigQuery. Historically, Amazon started AWS to monetize the additional (unused) infrastructure built to host Amazon's popular web store. Similarly, Google began competing with Amazon by offering a access to Dremel which, in turn, resulted in Google's Cloud Platform and the release of BigQuery.

Appendix

Appendix, INL 880: Capstone Product Worksheet (Submitted Separately)

References

Amazon Web Services. (2015) AWS Cloud Security. Amazon Web Services. Retrieved on June 7, 2015 from: <http://aws.amazon.com/security/>

Amazon Web Services. (2015) AWS Compliance. Amazon Web Services. Retrieved on June 23, 2015 from: <http://aws.amazon.com/compliance/>

Amazon Web Services. (2015) Amazon Redshift Pricing. Amazon Web Services. Retrieved on May 29, 2015 from: <http://aws.amazon.com/redshift/pricing/>

Amazon Web Services. (2015) Amazon Web Services Simple Monthly Calculator. Amazon Web Services. Retrieved on May 29, 2015 from: <https://calculator.s3.amazonaws.com/index.html>

Amazon Web Services (2015) AWS Support Pricing. Amazon Web Services. Retrieved on June 8, 2015 from: <https://aws.amazon.com/premiumsupport/pricing/>

Amazon Web Services. (2015) Amazon Redshift FAQs. Retrieved on July 15, 2015 from:
<http://aws.amazon.com/redshift/faqs/>

Ansari, Z., Azeem, M. F., Babu, A. V., & Ahmed, W. (2011). A fuzzy clustering based approach for mining usage profiles from web log data. International Journal of Computer Science and Information Security, 9(6), 70-79. Retrieved from
<http://search.proquest.com/docview/922726232?accountid=11999>

Atlassian Confluence Open Source Project. (October. 2014) Language Manual Joins. Apache Software Foundation. Retrieved on June 12, 2105 from:
<https://cwiki.apache.org/confluence/display/Hive/LanguageManual+Joins>

Apache Hive TM. (2014) Getting Started. The Apache Software Foundation. Retrieved on June 21, 2015 from: <https://hive.apache.org/>

awsdocumentation. (2015) Amazon Redshift Security Overview. Amazon Web Services. Retrieved on June 7, 2015 from: http://docs.aws.amazon.com/redshift/latest/dg/c_security-overview.html

awsdocumentation. (2015) Amazon Redshift: Database Developer Guide. Amazon Web Services. Retrieved on May 31, 2015 from:
<http://docs.aws.amazon.com/redshift/latest/dg/welcome.html>

awsdocumentation. (2015) Amazon Redshift: About Clusters and Nodes. Amazon Web Services. Retrieved on June 7, 2015 from : <http://docs.aws.amazon.com/redshift/latest/mgmt/working-with-clusters.html#rs-about-clusters-and-nodes>

awsdocumentation. (2015) Amazon Redshift: Analyzing the Query Plan. Amazon Web Services. Retrieved on June 7, 2015 from: <http://docs.aws.amazon.com/redshift/latest/dg/c-analyzing-the-query-plan.html>

awsdocumentation. (2015) Use a Multi-Row Insert. Amazon Web Services. Retrieved on June 7, 2015 from: http://docs.aws.amazon.com/redshift/latest/dg/c_best-practices-multi-row-inserts.html

Bauer, S. (2013). Getting started with amazon redshift. GB: Packt Publishing.

BigQuery. (2015) Third-Party Tools and Services. Google Inc. Retrieved on July 15, 2015 from: <https://cloud.google.com/bigquery/third-party-tools>

Borthakur, Dhruba (December 10, 2010) Looking at the code behind our three uses of Apache Hadoop. Facebook Inc. Retrieved on June 12, 2015 from: <https://www.facebook.com/notes/facebook-engineering/looking-at-the-code-behind-our-three-uses-of-apache-hadoop/468211193919>

Columbus, Louis. (Nov 22, 2014) Cloud Computing Adoption Continues Accelerating In The Enterprise. Forbes. Retrieved on May 24, 2015 from:

<http://www.forbes.com/sites/louiscolumbus/2014/11/22/cloud-computing-adoption-continues-accelerating-in-the-enterprise/>

Davison, Dean. (2014) The total economic impact of Google Cloud Platform. Forrester.

Retrieved on June 21, 2015 from: <https://cloud.google.com/why-google/analyst-reports/2014-03-Google-Cloud-Platform-TEI-Study.pdf>

Dean, Jeffrey; Ghemawat, Sanjay (2004). MapReduce: Simplified Data Processing on Large Clusters. Google. Retrieved on May 25, 2015 from:

<http://research.google.com/archive/mapreduce.html>

Ghemwat, Sanjay; Gobioff, Howard; Leung, Shun-Tak. (2003) The Google File System. Google.

Retrieved on May 25, 2015 from: <http://research.google.com/archive/gfs.html>

Google Cloud Platform (2015) Regions and Zones. Google Inc. Retrieved on June 8, 2015 from:

<https://cloud.google.com/compute/docs/zones>

Google Cloud Platform. (March. 2015) Google Security Whitepaper. Google Inc. Retrieved on June 12, 2015 from: <https://cloud.google.com/security/whitepaper>

Google Cloud Platform (2015) BigQuery: Access Control. Google Inc. Retrieved on July 14, 2015 from: <https://cloud.google.com/bigquery/access-control>

Google Cloud Platform (2015) BigQuery: Query Reference. Google Inc. Retrieved on June 8, 2015 from: <https://cloud.google.com/bigquery/query-reference>

Google Inc. (2015) “Why Google Cloud Platform” Retrieved on June 21, 2015 from: <https://cloud.google.com/why-google/>

Hertzfeld, E. (2015). From CapEx to OpEx. *Hotel Management*, 230(5), 25.

Hadoop. (2014) Welcome to Apache Hadoop! What is Apache Hadoop? The Apache Software Foundation. Retrieved on June 12, 2015 from: <https://hadoop.apache.org/>

IDG_Enterprise. (2014) IDG Enterprise Cloud Computing Infographic 2014. Scribd. Retrieved on May 24, 2015 from: <https://www.scribd.com/doc/246109036/IDG-Enterprise-Cloud-Computing-Infographic-2014>

jaql (2014) JaqlOverview; Jaql: A JSON Query Language. Google Inc. Retrieved on June 12, 2015 from: <https://code.google.com/p/jaql/wiki/JaqlOverview>

Keyser, Chris (May 11, 2015) Optimizing for Star Schemas and Interleaved Sorting on Amazon Redshift. Amazon Web Services. Retrieved on June 7, 2015 from:

7/20/2015 INL-880 - Capstone Proposal: McGinley & Etter -Final Draft

<https://blogs.aws.amazon.com/bigdata/post/Tx1WZP38ERPGK5K/Optimizing-for-Star-Schemas-and-Interleaved-Sorting-on-Amazon-Redshift>

Krause, R. (2011, Aug 02). Rackspace uses open source vs. amazon in cloud field rivalry fast emerging earnings out thursday as cloud vendor building support for OpenStack. Investor's Business Daily Retrieved from

<http://search.proquest.com/docview/928460241?accountid=11999>

Lampitt, A. (2012). Amazon's redshift brings cheaper big data analytics to small businesses. InfoWorld.Com, Retrieved on June 1, 2015 from

<http://search.proquest.com/docview/1220492610?accountid=11999>

Leoncini, R., Rentocchini, F., & Marzetti, G. V. (2011). COEXISTENCE AND MARKET TIPPING IN A DIFFUSION MODEL OF OPEN SOURCE VS. PROPRIETARY SOFTWARE (1). *Revue d'Économie Industrielle*, (136), 141-168. Retrieved from

<http://search.proquest.com/docview/1586117389?accountid=11999>

Leong, Lydia; Toombs, Douglas; Gill, Bob; Petri, Gregor; Haynes, Tiny. (May 28, 2014) ID: G00261698: Magic Quadrant for Cloud Infrastructure as a Service. Gartner Inc. Retrieved on May 24, 2015 from: <http://www.gartner.com/technology/reprints.do?id=1-1UKQQCY&ct=140528&st=sb>

Liebowitz, Jay. (Jun 13, 2013) Big Data and Business Analytics. CRC Press

Lohr, Steve. (April 10, 2012) I.B.M Aims to Sharply Simplify Corporate Data Center

Technology. The New York Times. Retrieved on May 24, 2015 from:

http://www.nytimes.com/2012/04/11/technology/ibm-aims-to-sharply-simplify-corporate-data-center-technology.html?_r=0

Marek, K. (2011). Web analytics overview. Library Technology Reports, 47(5), 5-10,2.

Retrieved from <http://search.proquest.com/docview/880109210?accountid=11999>

Melnik, Sergey; Gubarev, Andrey; Long, Jing Jing; Romer, Geoffrey; Shivakumar, Shiva;

Tolton, Matt; Vassilakis, Theo. (2010) Dremel: Interactive Analysis of Web-Scale Datasets.

Google Inc. Retrieved on June 8, 2015 from:

https://www.google.com/search?q=Dremel%3A+Interactive+Analysis+of+Web-Scale+Datasets&oq=Dremel%3A+Interactive+Analysis+of+Web-Scale+Datasets&aqs=chrome..69i57j69i58.280j0j7&sourceid=chrome&es_sm=91&ie=UTF-8

Minelli, M., Chambers, M., & Dhiraj, A. (2012). Wiley CIO, Volume 578 : Big Data, Big Analytics : Emerging Business Intelligence and Analytic Trends for Today's Businesses.

Somerset, NJ, USA: John Wiley & Sons. Retrieved from <http://www.ebrary.com>

Moore, Dennis. (June 25, 2011) The real (potential) impact of SAP HANA. Enterprise Irregulars. Retrieved on June 7, 2015 from: <https://www.enterpriseirregulars.com/39209/the-real-potential-impact-of-sap-hana/>

Mone, G. (2013). Beyond hadoop. New York: ACM. doi:10.1145/2398356.2398364

MongoDB (2015) A Total Cost of Ownership Comparison of MongoDB & Oracle. Retrieved on May 31, 2015 from <https://www.mongodb.com/collateral/total-cost-ownership-comparison-mongodb-oracle>

Oracle (2000). Oracle8i Integration Server Overview Release 3 (8.1.7) Part Number A83729-01: Directory Services (LDAP). Oracle. Retrieved on June 8, 2015 from: http://docs.oracle.com/cd/A87860_01/doc/ois.817/a83729/adois09.htm

Peters, Brad. (Feb, 28, 2013) And the Winner is, Who?! Forbes. Retrieved on May 31, 2015 from : <http://www.forbes.com/sites/bradpeters/2013/02/28/and-the-winner-is-who/>

Progressive Digital Media Technology News (Dec 17, 2014). 10 hadoop predictions for 2015. Progressive Digital Media Technology News. Retrieved on June 5, 2015 from: <http://search.proquest.com/docview/1637485778?accountid=1199>

Raab, D. (2007). How to judge a columnar database. DM Review, 17(12), 33. Retrieved on June 3, 2015 from <http://search.proquest.com/docview/214670809?accountid=11999>
7/20/2015 INL-880 - Capstone Proposal: McGinley & Etter -Final Draft

Rajaraman, V. (2014). Cloud computing. Resonance: Journal Of Science Education, 19(3), 242-258. doi:10.1007/s12045-014-0030-1

Rouse, Margaret. (2006) Key Performance Indicator (KPI). TechTarget. Retrieved on June 23, 2016 from: <http://searchcrm.techtarget.com/definition/key-performance-indicator>

Russom, Philip. (2013) Integrating Hadoop Into Business Intelligence and Data Warehousing. TDWI Research. Retrieved on July 15, 2015 from:
<https://www.cloudera.com/content/dam/cloudera/Resources/PDF/TDWI%20Best%20Practices%20report%20-%20Hadoop%20for%20BI%20and%20DW%20-%20April%202013.pdf>

Sato, Kazunori (2012) An Inside Look at Google BigQuery. Google Inc. Retrieved on June 7, 2015 from <https://cloud.google.com/files/BigQueryTechnicalWP.pdf>

Scott, Peter (Feb, 20, 2014) Thoughts on Using Amazon Redshift as a Replacement for an Oracle Data Warehouse. Rittman Mead. Retrieved on May 31, 2015 from:
<http://www.rittmanmead.com/2014/02/thoughts-on-using-amazon-redshift-as-a-replacement-for-an-oracle-data-warehouse/>

Spry. (2013) A list of Subtle Difference Between HiveQL and SQL. Spry. Retrieved on June 8, 2015 from: <http://spryinc.com/blog/list-subtle-differences-between-hiveql-and-sql>

Turkington, G. (2013). Hadoop Beginner's Guide. Olton, Birmingham, GBR: Packt Publishing.

Retrieved from <http://www.ebrary.com>

Vambenepe, William (April 16, 2015) Big Data, the cloud way. Google Inc. Retrieved on June 7, 2015 from: <http://googlecloudplatform.blogspot.com/2015/04/big-data-cloud-way.html>

Ward, Miles (January 28, 2015) Understanding Cloud Pricing. Google Inc. Retrieved on June 5, 2015, from: <http://googlecloudplatform.blogspot.com/2015/01/understanding-cloud-pricing.html>

Ward, Miles (June 19, 2015) Understanding Cloud Pricing Part 3 - Data Warehouses. Google Inc., Retrieved on June 19, 2015 from:

Ward, Miles (June #, 2015) Understanding Cloud Pricing Part 3.2 - More Data Warehouses. Google Inc., Retrieved on June #, 2015 from:

Watson, L. A., & Mishler, Chris, CMA,C.I.A., C.I.S.A. (2014). From on-premise applications to the cloud. Strategic Finance,96(2), 80-81. Retrieved from <http://search.proquest.com/docview/1552717174?accountid=11999>

What's A Byte? (2015) "Megabytes, Gigabytes, Terabytes... What Are They?" Retrieved on July 15, 2015 from: <http://www.whatsabyte.com/>